

Working Memory

Bradley R. Postle and Klaus Oberauer

To appear in: M.J. Kahana and A.D. Wagner (Eds.) *The Oxford Handbook of Human Memory*.
Oxford University Press (Oxford, U.K.).

Bradley R. Postle, PhD
Departments of Psychology and Psychiatry
University of Wisconsin–Madison
Madison, WI. 53706
USA
postle@wisc.edu

Klaus Oberauer, PhD
Department of Psychology
Cognition Psychology Unit, University of Zurich
Zurich
Switzerland
k.oberauer@psychologie.uzh.ch

1. Introduction

Working memory refers to the ability to hold information in an accessible state – in the absence of relevant sensory input – to transform it when necessary, and to use it to guide behavior in a flexible, context-dependent manner. Individual differences in working memory ability are relatively stable and trait-like, and they predict an impressive array of laboratory measures and real-world outcomes, from general fluid intelligence¹⁻⁵ to reading⁶ to scholastic achievement⁷. It can be demonstrated with any modality of sensory information, alone or in combination, as well as for most domains of cognition. Because working memory is understood to be a necessary elemental contributor to many aspects of high-level cognition – such as cognitive control, problem solving, and planning e.g.,^{8,9,10} – and because its impairment is characteristic of many neurological and psychiatric syndromes – including attention deficit hyperactivity disorder (ADHD), Alzheimer’s disease (AD), Parkinson’s disease (PD), major depressive disorder (MDD), and schizophrenia e.g.,^{11,12} – it is the focus of intensive study within several domains of psychology, neuroscience, and medicine. Across these disciplines, working memory has been studied most intensively with tasks requiring visual, auditory, and linguistic processing, and these will be the focus of this review.

The formal study of the ability to mentally retain short lists of verbal material dates back at least to the time of Ebbinghaus¹³ and James¹⁴. The idea of working memory as a key element of cognition, however, emerged during the cognitive revolution, as psychologists began to explicitly consider cognition from an information-processing perspective. Beginning in the 1950s, computational models of human problem solving incorporated a “working memory” that served a function similar to that of random access memory (RAM) in the architectures of computing machines^{15,16}. This, in turn, influenced the thinking of Miller, Galanter, and Pribram (1960)¹⁷ in their articulation of an alternative to the behaviorist program for explaining higher-level cognition. In their conception, even the most elemental processing of an input involved its comparison against an internal model, the outcome of which would determine an organism’s response to that input. This necessitated the incorporation of feedback, and meant that the analysis of even the simplest of sensory-motor events needed to incorporate principles from information theory¹⁸, and an appreciation that all cognitive processing involves the implementation of control. When considering higher levels of cognition (e.g., planning, decision making, communicating), Miller et al. (1960)¹⁷ characterized the propositional units of cognition as “Plans,” and asserted the following about the execution of a Plan:

“... something important ... happen[s] to a plan when the decision is made to execute it. It is taken out of dead storage and placed in control of a segment of our information-processing capacity. It is brought into the focus of attention, and as we begin to execute it we take on a number of menial but necessary tasks having to do with gathering data and remembering how far in the Plan we have progressed at any given instant, etc. Usually the Plan will be competing with other Plans also in the process of execution, and considerable thought may be required in order to use the behavioral stream for advancing several Plans simultaneously. The parts of a Plan that is being executed have special access to consciousness and special ways of being remembered that are necessary for coordinating parts of different Plans ... When we have decided to execute some particular Plan, it is probably put into some special state or place where it can be remembered while it is being executed... Without committing ourselves to a specific machinery, therefore, we should

like to speak of the memory we use for the execution of our Plans as a kind of quick-access, ‘working memory.’ (p. 65)”

This passage invokes several concepts that remain highly relevant for contemporary models of and debates about working memory, and we will consider many of them over the course of this chapter: the activation of information from long-term memory (LTM); the focus of attention; the distinction between the rules and/or goals that are guiding behavior (i.e., the “Plans”) versus the situation-specific information whose influence on behavior is determined by those rules (i.e., the “gathering of data”); competition and interference between mental representations and action plans; the relation between working memory and conscious awareness; and the question of whether holding information in working memory entails putting it into a “special state” or a “special place.” Indeed, elaboration on two of these points will allow us to highlight two themes that will be relevant for each of the topics that will be addressed in this chapter.

First, “Plans” versus “the gathering of data.” This highlights the fact that the control of behavior is often guided by hierarchically organized rules, and that the rules governing behavior in a particular situation may be processed differently than is the information (the “data”) being held in working memory. When driving a car on an unfamiliar road, for example, the information conveyed by the just-passed road sign -- that the second exit off the roundabout leads to your destination -- will typically guide your immediate behavior. If your car is also low on petrol, however, and you see that there is a petrol station at the third exit off this roundabout, the information in working memory will influence your behavior differently: you will use this information to prompt yourself to make note of the distinctive landmarks at the second exit as you drive past it on your way to the petrol station. To translate this distinction to the laboratory, in tests of digit span, for example, it is useful to distinguish between the content of working memory on any given trial, which is the series of digits spoken by the experimenter, and the rule governing behavior, which would be whether the subject is to recall the digits in the order in which they were presented (“forward digit span”) or in the reverse order (“backward digit span”). The analogous distinction can also be made for laboratory tests of nonhuman animals, where the content of a trial may be the location briefly cued on a screen, and the rule whether the subsequent delayed saccade is to be made to the cued location (a pro-saccade) or to a location 180° opposite of the cue (an anti-saccade); or where the content may be the sample object presented at the beginning of the trial, and the rule whether the subject is to select that sample when it is re-presented in a test array of two objects (“delayed match-to-sample”), or to select the novel object (“delayed non-match-to-sample”). Working memory for rules versus for content can be dissociated neurally, and the two may differ, in some circumstances, in terms of their access to conscious awareness. Furthermore, as we shall see, these and other considerations have led to the proposal of a fundamental distinction between a procedural working memory versus a declarative working memory.

The ‘state or place’ question also merits additional consideration in this introductory section, because it gets to a fundamental question of how to best situate working memory within the broader context of cognition, as well as within the neural systems that underlie cognition. On one hand, if one assumes that information being held in working memory has been “put into some special place,” one is assuming that there is a dedicated mechanism, with one or more identifiable sub-system(s) of the mind/brain, that serves this function. Many contemporary models take such a memory-systems perspective, as exemplified by models

positing working memory buffers responsible for the domain-specific storage of the contents of working memory, and by the localization of this buffering function to sustained, elevated activity in the prefrontal cortex (PFC). Alternatively, we could assume that holding information in working memory can be accomplished by it having been “put into some special state”. This special state can be understood as a property of the system that is specialized for representing the information in question, such as the systems that support sensory perception, skeleto- or oculomotor control, language, or semantic memory. From this state-dependent perspective, one would expect the content of working memory to be retained via a transient state change (perhaps sustained activation, perhaps modified synaptic weighting) of the same representational systems that process this information in contexts that do not make overt demands on working memory, such as perceiving, carrying out an action, or thinking about facts about the world. Whether we need to assume any mechanisms specific to working memory, or whether working memory can be fully explained as emerging from the operation of systems of the mind/brain that evolved for other functions, is a topic of ongoing debate to which we will return periodically in this chapter. For now, we define working memory by its function, leaving open whether this function is fulfilled by a dedicated system or by the cooperation of other cognitive/neural systems.

This chapter begins with an overview of the functional requirements for working memory, and some of its cardinal properties. Next we will briefly summarize a few theoretical models that exemplify current conceptualizations from memory-systems and state-dependent perspectives. Section 4 will dig deeper into three questions of considerable interest in contemporary working memory research, and the final section will provide an overview of the neural bases of working memory functions.

2. Functional Requirements for Working Memory

As sketched out in the Introduction, we define working memory by its function, which is to guide our current thoughts and actions with temporarily selected representations -- in the terminology of Miller, Galanter, and Pribram (1960)¹⁷, to carry out “Plans” by processing the recently gathered “data”. To fulfill this function efficiently, working memory should have the following characteristics¹⁹: (1) There should be a medium for rapidly building and maintaining temporary bindings between representations, so that the existing representational units (chunks) in LTM can be combined into new structures. For example, when one needs to briefly remember a novel telephone number, or an array of colored squares, the tokens themselves (i.e., the digits or the colors) are already familiar and represented in LTM – what determines success or failure in these situations is remembering the order of these digits that corresponds to the phone number, or the location of each color that is specific to this array. We will refer to these new structures that control our cognitive activity as the representations “in working memory” without implying that they are maintained in a dedicated buffer. (2) There should be mechanisms for manipulating these structural representations. These include (a) an attentional mechanism for selective access to those contents of working memory that need to be manipulated next, and (b) a mechanism for maintaining procedural representations (i.e., plans, goals, rules) in working memory that control how the declarative contents of working memory are to be manipulated. (3) Working memory needs to hold the information most relevant for the cognitive system's current goals at any point in time. This entails two conflicting demands: On the one hand, relevant information must be protected against interference from other, irrelevant information. This can be accomplished by shielding working memory against input from perception and from LTM. On the other hand, working memory contents need to be rapidly updated: New relevant contents need to be encoded

quickly, and old, no longer relevant contents removed quickly from working memory. Together, these requirements pose a stability-flexibility dilemma: Contents need to be stably maintained and shielded from interference as long as they are relevant, but rapidly removed and replaced when they become irrelevant, or when other information takes higher priority. To meet both demands, working memory needs to shift flexibly between a maintenance mode and an updating mode^{20,21}. We next discuss these three requirements in turn and review evidence speaking to them.

2.1 Temporary Bindings

Language processing, planning, reasoning, and problem solving all involve the construction of new representations by combining known elements in novel ways. Sentences are formed by combining words in a novel order; an action plan assembles familiar steps into a sequence; reasoning about a mechanical device involves construction of a mental model from known mechanical elements and forces. Assembling new structural representations from known elements requires a mechanism for the rapid formation of temporary bindings. Theories of reasoning and language processing have long acknowledged the key role of bindings between content elements and their places or roles in a structure^{22,23}. Representing the meaning of a sentence involves binding the concepts referred to by the content words to their roles in a proposition. For instance, understanding "The dog chases the cat" involves binding the concept DOG to the agent role, the concept CAT to the object role, and CHASE to the action role in a proposition. Constructing an action plan involves binding each action to its ordinal position in the planned sequence of steps. The elements of a mental model of a mechanical device are bound to their spatial locations and to their roles in the causal chain or network governing the device.

Experimental studies of working memory usually ask participants to briefly hold in mind comparatively simple structural representations, such as the serial order of words in a list, or the spatial arrangement of colors in an array. Representing these memory sets in working memory involves binding each element to its location in a mental space (i.e., to their position in a list, or their spatial location in an array). Many computational models of working memory make this binding mechanism explicit (see *section 3.3*). For instance, the most successful models of serial recall of lists share the assumption that list items are bound to their positions on a dimension of psychological time²⁴⁻²⁶ or to their ordinal positions in the event sequence^{27,28}. Recent models of visual working memory also incorporate bindings between visual objects and their spatial locations^{29,30}.

Experiments have shown that failures of bindings are responsible for a large proportion of errors in working memory tasks. For instance, when people try to recall a list in order, they often report the list elements in the wrong order. These order errors are most often confusions between elements in positions close to each other, a tendency referred to as locality constraint³¹. Binding failures also account for a substantial proportion of errors in tests of visual working memory^{32,33}. Similar to the locality constraint in time for lists, a locality constraint in space has been observed for confusions between stimuli in visual arrays: Elements are most likely to be confused with close neighbors^{34,35}. The locality constraint shows that binding failures arise in part because the representations of each element's position in time or space is imprecise.

2.2. The "Working" of Working Memory: Mental Manipulation

Mental operations such as language processing, planning, reasoning, and problem solving involve not only constructing structural representations in working memory but also working with them. We next discuss two mechanisms that enable mental work: Selective access to subsets of the contents of working memory, and representations of task goals and rules in procedural working memory that control the mental operations on the declarative contents of working memory.

2.2.1. Selective Attention to Elements in Working Memory

Working with the contents of working memory typically means to operate on individual elements, or subsets of elements, of the structure currently held in working memory. For instance, after planning a short sentence to be spoken, the person will want to say each word in order. To do this, they need to select one word at a time from the ordered set of words in working memory. Selective access to elements in a memory set is a form of attention directed to working memory representations. Several lines of research have confirmed that people can direct attention to elements within a memory set, thereby temporarily prioritizing it without forgetting the other elements in the set. For instance, after encoding a list of items, the last-encoded item remains in a state of particularly fast accessibility for a few seconds³⁶⁻³⁸. This advantage for the last-encoded item can also be shown at each step during encoding of a memory list when memory is probed in between presentation of one item and the next³⁹. More generally, the last-retrieved or last-updated element in a memory set remains in a state of fast accessibility^{40,41}. These findings can be explained by the assumption that the last-encoded or last-used item remains for a while in a focus of attention within working memory, understood as a qualitatively special state of being selected for processing. However, these findings can also be interpreted as reflecting a steep recency gradient of memory strength⁴².

Evidence for an attentional selection mechanism in working memory also comes from the retrodictive cuing (retro-cue) effect^{43,44}. Retro-cues have mostly been studied in the context of visual-working memory tests. After encoding an array of visual objects, a cue directs attention to the one item that will most likely be tested. If the cue validly identifies the item that is tested, responses become faster and more accurate; in case of an invalid cue, performance is impaired relative to a no-cue baseline condition. These retro-cue effects are observed when the cue is presented one second or more after offset of the memory display, ruling out the possibility that the cue taps into sensory memory. The available evidence (reviewed in⁴⁵) suggests that several mechanisms are jointly responsible for the retro-cue effects: Attending to an element in working memory strengthens the binding of that element to its location in the array; protects that element against interference from further visual input; and sometimes triggers removal of the other, not-cued elements from working memory.

It is tempting to think that the two sets of findings – facilitated access to the last-used item, and the retro-cue effect – point to the same attentional mechanism in working memory. This is doubtful, however, in light of experiments showing that the two effects combine additively: A retro-cue to the last-presented item of a list boosts access to that item as much as a retro-cue to earlier-presented list items⁴⁶. If a retro-cue brings the cued item into the focus of attention, and the last-presented item is already in the focus of attention, then a retro-cue to the last-presented item could not add anything to its privileged status. Therefore, it appears more likely that the heightened accessibility of the last-used item reflects a recency gradient of memory strength, and the retro-cue effect arises from attentional selection of the cued item.

2.2.2: Procedural Working Memory

As foreshadowed by Miller, Galanter, and Pribram (1960)¹⁷, working memory needs to hold two kinds of representations: The contents to be manipulated (the "data"); and representations that control how these contents are manipulated (the "Plans" according to Miller et al. (1960), or "task sets" in more recent terminology). Borrowing from theories of long-term knowledge, we refer to the former as "declarative" and the latter as "procedural" representations in working memory⁴⁷.

Procedural representations in working memory can be described as "if-then" rules, which link a condition to an action (in this regard they are like productions in production-system models of the mind, e.g., ⁴⁸). When the condition is met by the current declarative content of working memory (in particular, the content selected into the focus of attention) then the action – which could be a mental manipulation of the declarative working memory contents, or a physical action guided by these contents – is carried out. There are many such rules in a person's knowledge repertoire, and often these rules can be in conflict with each other. For instance, a person commuting between the UK and continental Europe has learned the rule "if you drive, stay on the right side of the street" as well as the rule "if you drive, stay on the left side of the street". Which rule applies depends on the context (in this example, on the country the driver is in). People can rapidly switch between alternative (potentially conflicting) rules to be applied to the same situation, a feat that has been extensively studied in the literature on task switching⁴⁹⁻⁵¹. To do this, the mind needs to select at any point in time one rule or task set that is to govern mental and physical action, at the exclusion of other potentially competing task sets that could be applied to the same situation. The task set selected for this purpose is the procedural representation in working memory.

A task set can be established in working memory by retrieving it from LTM⁵². It can also be created in working memory to implement a new instruction, such as "if you see a picture of a four-legged animal, press the left button, and if you see a two-legged animal, press the right button". Humans can implement arbitrary instructions like these as procedural representations in working memory without practice. This is shown by the fact that when people receive a new instruction mapping stimulus categories to responses, the instructed rules interfere with an ongoing task even before they have ever been carried out⁵³. This observation can be explained by the assumption that instructed rules are established as procedural representations in working memory, which operate as a "prepared reflex": Whenever their condition is met by a stimulus that a person attends to (or another declarative representation in working memory), the action bound to it is carried out automatically.

2.3. Meeting the Stability-Flexibility Dilemma: Gated Encoding and Updating

To control our cognitive processes such that they serve our current goals, working memory needs to hold available the information most relevant for these goals at a given time. This sometimes means to maintain information for a while even in the absence of supporting input from perception or LTM, for instance when we plan to come back to a topic while the conversation drifts to another topic. On other occasions it means to rapidly discard information that is no longer needed, such as a phone number you no longer need to dial because someone else in the room got to their phone first. Hence, working memory needs to meet opposing demands: To maintain relevant information it is best to close the gate to any further input so that the current content is shielded from interference. To seamlessly update working memory it is necessary to open the gate to new input, and to rapidly remove the current information. Sometimes both demands arise at the same time because part of the

current contents of working memory need to be updated while others need to be retained. For instance, imagine your colleague just told you her new phone number: 326 74 24, but then corrects herself: "actually, no, the last two digits are 59". You need to first build a seven-digit list in working memory, and then selectively replace the last two digits while keeping the first part of the list. To meet these demands, working memory needs mechanisms for gating its input from both perception and LTM, and mechanisms for efficiently removing no-longer-relevant information. We next review evidence for both mechanisms in turn.

2.3.1. Gated Encoding

Young adults are very good – though not perfect – at limiting encoding into working memory to those aspects of the perceived environment that they deem relevant for the current task. In the cognitive control literature, this is sometimes referred to as “input gating”. For instance, when presented with a list of words, and instructed to remember every second word, they can repeat the relevant words nearly as well as if only the relevant words had been presented⁵⁴. Keeping irrelevant verbal information out of working memory is harder when it is spoken: Working memory maintenance of verbal lists is substantially impaired by concurrent irrelevant speech, and also non-speech sound streams with high variability^{55,56}. However, this impairment does not appear to arise from the irrelevant sound interfering with working memory contents because the effect is independent of the similarity between to-be-remembered and to-be-ignored stimuli⁵⁷. One explanation for the irrelevant-sound effect is that the irrelevant stream, by being a sequence of events itself, interferes with the mechanism of maintaining the serial order of the memory items^{58,59}. An alternative account is that the irrelevant sound captures attention, thereby disrupting attention-dependent processes of encoding and maintenance of the relevant material^{60,61}.

When presented with an array of oriented bars or triangles – some red, some blue – and instructed to remember only the red ones, their performance is only slightly impaired by the presence of blue stimuli compared to arrays containing only red stimuli⁶². That said, an EEG-based indicator of working memory load, the contralateral delay activity (CDA, see *section 5.3.1*), shows that irrelevant stimuli are encoded into working memory to some extent, and more so in people with smaller estimates of visual working memory capacity⁶³.

Keeping information out of working memory is harder when it is relevant at least temporarily: When during maintenance of a memory set some distractor stimuli need to be processed (e.g., reading words aloud or making a judgment on them) but not remembered, these distractors nevertheless are encoded into working memory²⁸. Their strength of encoding, however, can be reduced to about half of that of the memory items⁶⁴.

Retrieval cues can bring information in LTM to mind automatically. Thus, the contents of working memory also need to be shielded against irrelevant, potentially interfering information from LTM. At the same time, however, working memory should be open to relevant or potentially helpful information from LTM. Ideally, a flexible gate should admit information from LTM into working memory to the extent that this information is helpful rather than interfering. In fact, there is a wealth of evidence for facilitating influences of knowledge in LTM on maintenance in working memory. Memory sets matching known units in semantic LTM (such as the letter sequence “PDF”) are remembered better than memory sets not matching any knowledge (e.g., the sequence “FPD”)⁶⁵; lists of words are remembered better than lists of pseudowords⁶⁶, and memory lists repeated several times

across an experiment are recalled better from trial to trial⁶⁷. In contrast, there is little, if any, evidence for LTM contents interfering with maintenance of information in working memory (for a review see ⁶⁸). In experiments directly comparing facilitating and interfering effects of LTM knowledge on maintenance in working memory, there is evidence for facilitation, but against interference, as predicted from the assumption of a flexible mechanism that uses LTM information only in situations when it is helpful rather than harmful⁶⁸.

2.3.2 Updating of Working Memory

Because the contents of working memory are the contents of our current thoughts, working memory needs to be updated at the speed at which our thoughts progress – that is, several times per second. Replacing the entire working memory content by a new memory set is a fast process; partial updating of some elements while maintaining others is considerably slower⁶⁹, reflecting the challenge of balancing stability and flexibility at the same time.

A detailed analysis of response-times across several conditions revealed how this challenge is met^{70,71}: When participants hold a list of letters in working memory, and are asked to replace a subset of them by new letters presented on the screen, they scan the letters in the habitual reading direction associated with the material – left-to-right for Latin letters, right-to-left for Hebrew letters. At each step a decision is made whether the currently focused element is to be maintained or to be removed. A large part of the time demand of selective updating is due to the time cost of switching between maintenance and substitution. Subsequent experiments using another working memory-updating paradigm confirmed that switching between a maintenance mode and an updating mode of working memory incurs a substantial switch cost⁷².

2.4. Summary

We have argued that the function of working memory is not primarily to remember information but to hold it available to control information processing, as a "Plan" in the sense of Miller, Galanter, and Pribram (1960)¹⁷. This function entails several requirements: Working memory needs a mechanism for forming, maintaining, and flexibly updating bindings; it needs a mechanism for selectively accessing subsets of its contents; it must be able to hold procedural as well as declarative representations, and its contents need to be shielded to some extent against influences from perception and long-term memory while being open to relevant input from both channels. Theories and computational models of working memory – reviewed next – reflect some of these requirements, although we are still far from understanding the mechanisms working together to meet them.

3. Theories of working memory

A thorough review of theories and models of working memory is outside the scope of this chapter^a. Here we will introduce just a small number of theoretical frameworks that will be useful for contextualizing the current state of the literature, followed by a brief survey of computational models of working memory. In particular, we summarize key tenets of

^a One collection that assembles influential models from the turn of the century is 73. Miyake A, Shah P, eds. *Models of Working Memory*. Cambridge, U.K.: Cambridge University Press; 1999., and an authoritative update is expected to be published at around the same time as will be this volume 74. Logie R, Camos V, Cowan N, eds. *Working Memory: State of the Science*. Oxford University Press; in preparation.

currently influential theories that exemplify memory-systems and state-dependent perspectives.

3.1 Memory-systems models

Far-and-away the most influential memory-systems model, and arguably the single most influential model in the modern study of working memory, is the multiple-component model first proposed by Baddeley and Hitch⁷⁵ and subsequently updated on several occasions by Baddeley and collaborators (e.g.,^{8,76-79}). This model posits several domain-specific memory buffers – a phonological loop for verbalizable information, a visuospatial sketchpad for visuospatial information, and an episodic buffer for information retrieved from LTM – and a Central Executive responsible for the coordination of the operations of these buffers and for the manipulation of their contents.

At a finer grain of detail, the phonological loop is comprised of a phonological store that provides the storage function, and an articulatory loop responsible for rehearsal. Verbalizable information that is presented acoustically has obligatory access to the phonological store, as does visually presented information (e.g., written words, letters, or digits) after it is automatically recoded into a phonological code. Retention in the phonological store is subject to a decay factor, degrading information to an irrecoverable state within 2 sec⁸. The effects of decay can be counteracted by the periodic refreshing effects of rehearsal in the articulatory loop, the rate of which corresponds to the rate of overt articulation. Thus, the capacity of verbal working memory (classically understood to be 7 +/- 2 items for material processed in English⁶⁵) is explained by this model as the number of items that can be rehearsed within a 2-sec span. A concrete demonstration of this is the fact that the same (bilingual) individual will have a larger digit span when remembering digits in English than when remembering them in Welsh, because overt articulation in the latter language is slower, a factor assumed to also influence covert rehearsal⁸⁰. The operation of the articulatory loop can be blocked by concurrent overt articulation, an intervention that dramatically decreases verbal working memory span (e.g.⁸¹), and that also blocks access of visually presented verbalizable information to the phonological store⁸.

The architecture of the visuospatial sketchpad has been influenced by research on visual processing. The functional distinction between “what” an object is and “where” it is located⁸² is mirrored in the fractionation of the visuospatial sketchpad into a visual cache (for representing object features) and an inner scribe (for representing spatiotemporal information^{83,84}).

The episodic buffer was added to the multiple component model to provide a substrate for the linking of information across modalities into novel coherent representations (as, e.g., for a narrative sequence of events), and for buffering information retrieved from episodic and semantic LTM⁷⁶. The work of binding items in the episodic buffer to generate novel structural representations is assumed to require the Central Executive.

3.2 State-dependent models and frameworks

An early and influential articulation of a state-dependent model of working memory is Cowan's⁸⁵⁻⁸⁷ embedded-processes model, which explains working memory as arising from the temporary activation of preexisting knowledge structures within the cognitive system. Key to this model is the distinction between representations that are merely activated –

referred to as "activated LTM" (aLTM) – and representations in the Focus of Attention (FoA), a capacity-limited privileged state in which information must be held for its manipulation, and for access to awareness. The FoA is assumed to differ from aLTM in several regards: First, it has a discrete capacity limit of 3-4 chunks (see *section 4.1*). Second, whereas information in aLTM is subject to forgetting through decay and interference, the contents of the FoA are protected from these corrosive influences. Third, whereas aLTM only consists of already existing knowledge that is merely activated, individual chunks in the FoA can be bound together, thereby enabling the construction of new structural representations.

Another theoretical framework describing a state-dependent working memory has been proposed by Oberauer⁸⁸⁻⁹⁰. This framework distinguishes three states of information in working memory: activated LTM (similar to the corresponding concept in Cowan's model), the region of direct access, and the focus of attention. The region of direct access is a mechanism for creating and maintaining ad-hoc bindings, thereby forming new structural representations. Its capacity is limited by interference between bindings. The focus of attention is a selection mechanism that enables access to individual elements within the representational structure currently held in the region of direct access.

3.3 Computational models of working memory

On a more detailed level than the theoretical frameworks reviewed above, theoretical ideas about working memory have been expressed as computational models that describe the hypothetical mechanisms and processes underlying working memory as mathematical functions. Computational models have been proposed on several levels of granularity, ranging from abstract mathematical formulations^{91,92} to detailed simulations of neural networks⁹³⁻⁹⁵. There are two broad classes of working-memory models: Activation-based models and connection-based models.

Activation-based models of working memory^{95,96} build on the long-standing assumption in cognitive neuroscience that maintenance in of information in working memory relies on persistent firing of neural assemblies representing that information throughout the retention interval. The elementary computational units of these models are model neurons that sustain their activation through recurrent connections by which they re-activate themselves (directly or indirectly via other neurons), often accompanied by inhibitory connections between units enrolled in other representations. Information about an item in working memory (e.g., a word or a color) can be represented by a single unit or an assembly of units that re-activate each other through their connections. An influential model of this kind is the "bump attractor" model by Wei, Wang, and Wang⁹⁶. This model was built to explain the maintenance of simple visual features varying along a continuous dimension, such as colors and orientations, in working memory. The model architecture consists of a bank of model neurons, each of which has a tuning curve that describes its response to a stimulus. The tuning curve peaks at the "preferred" stimulus of a unit and gradually declines as stimuli become more dissimilar from the preferred stimulus (see Figure 2A). For circular feature dimensions such as orientations in 2-dimensional space, the neurons are arranged in a ring, ordered by the similarity of their preferred orientations. Neighboring units in this circular line-up are connected by excitatory and more distant units by inhibitory synaptic links. In this way, each stimulus creates a "bump" of activation centered on the model neuron that maximally responds to that stimulus (Figure 2B). When several orientations are to be encoded into working memory, several bumps are created simultaneously. Because bumps at different locations in the ring of neurons inhibit each other, the number of bumps that can be

upheld after stimulus offset is limited, thereby explaining the limited capacity of working memory.

One limitation of activation-based models is that they have no general mechanism for representing bindings (see *section 2.1*). Most working-memory tasks require maintaining information about arbitrary relations – for instance, when participants are asked to remember a random list of letters, they need to remember the relation of each letter to its ordinal position in the list. When they try to remember an array of colors, they need to represent the relation of each color to its location in the array. Activation-based models can only sustain activation of already existing representations but not represent new relations between them. Moreover, as reviewed in *section 5.3.3*, the assumption that maintenance in working memory relies on persistent neural activation has been questioned by recent developments in cognitive neuroscience. Rapid, temporary changes to the strength of synaptic connections between neurons have been discussed as an alternative^{97,98}.

Rapid changes of connection strengths form the backbone of a second class of computational models of working memory^{25,27,99}. They are often implemented as neural networks with two layers of units: One for representing the contents to be remembered (e.g., letters, colors), and the other the context that serves as the retrieval cue to access them when needed. At encoding, content representations are bound to their contexts through rapid changes of the connection weights between the two layers of units (see Figure 2C, D). For instance, letters of a list to be recalled in forward order are bound to their serial positions in the list, so that at test, the positions, activated in forward order, serve as retrieval cues to reproduce the letters. Similarly, visual objects presented in different locations are bound to their spatial locations; at test, the color or shape of an object can be retrieved when its location is given as a cue (or vice versa). Models of this kind have been applied successfully to two major experimental paradigms for studying working memory, serial recall of lists (for a more detailed treatment see the chapter of Hurlstone in this volume) and recall of visual objects in arrays.^{29,99,100}

4. Important Questions in Contemporary Working Memory Research

There are far too many topics of current research related to working memory for us to be able to cover in a single chapter. Here we limit ourselves to the three that we think of as the most important ones: (1) The nature of the capacity limitation of working memory; (2) the relation between working memory and LTM, and (3) the relation of working memory to attention and the control of thoughts and actions.

4.1 Capacity limits and the units of representation

The factors that account for the capacity limitations of working memory are the subject of intense study and debate. This question is taken up in greater detail in the chapter from this volume by Foster, Awh, and Vogel¹⁰¹, but here we will summarize some of the important points that make contact with other sections of this chapter.

Three classes of hypothesis have been proposed to explain why the capacity of working memory is limited¹⁰². These hypotheses have different implications for the units of measurement of capacity.

4.1.1. Time-based decay

The first hypothesis is that memory traces in working memory decay rapidly unless refreshed by some form of rehearsal, as assumed for the phonological loop component in the multiple-component model (see *section 3.1*). It follows from this hypothesis that capacity is best measured in units of time: The capacity of working memory is the duration for which information can be maintained in working memory without being rehearsed. If the speed of rehearsal for some class of materials is known, the capacity can be expressed as the amount of information that can be rehearsed within the time limit given by decay¹⁰³. This same principle also applies to the phonological loop component of the multiple-component model.

4.1.2. Limited resource hypotheses

The second class of hypothesis holds that working-memory capacity is constrained by a limited resource that has to be divided among all representations that need to be maintained simultaneously; in some models the same resource also needs to be shared with concurrent cognitive processes. The resource hypothesis comes in two flavors: Discrete and continuous resources.

4.1.2.1 Discrete resource

Discrete-resource models assume that the capacity of working memory is determined by a discrete number of place-holders (sometimes referred to as "slots"), each of which can hold one representational unit. This idea has been advocated by Cowan⁸¹, who reviewed a broad set of findings from various experimental paradigms challenging working memory and arrived at the conclusion that the capacity of working memory amounts to about 3-4 chunks on average in healthy young adults. Obviously, this notion implies that the unit of measurement for working-memory capacity is the number of chunks that can be maintained. This number can be estimated from performance in working-memory tasks through measurement models that incorporate assumptions about how people guess when tested for information that did not fit into a slot^{104,105}. The bump-attractor model described in *section 3.3* has been proposed as one mechanism for creating a slot-like capacity limit.

4.1.2.2 Continuous resource

Continuous-resource models assume that the resource can be continuously divided into arbitrarily small shares, and therefore there is no limit to the number of representational units that can be maintained in working memory¹⁰⁶. As the number of units increases, the share of resource that each of them receives decreases. The resource share assigned to a representation is monotonically related to the ability to retrieve it through a performance-resource function¹⁰⁷. In some models, the resource share determines the chance of retrieval¹⁰⁸, whereas in others it determines the precision of the retrieved information³². In these models, the capacity of working memory is measured as the total resource quantity available; it can be estimated from performance through measurement models that incorporate assumptions about the performance-resource function¹⁰⁹. One way in which a resource limit could arise in a neural network is by divisive normalization: The activation or firing rate of all neurons in a network is constrained so that their sum must not exceed a fixed maximum.⁹⁴

4.1.3 Interference

The third hypothesis is that interference between representations in working memory causes the capacity limit. This hypothesis has been fleshed out in computational models – briefly described in *section 3.3* – in which contents (such as words, digits, or visual objects)

are bound to contexts that serve as retrieval cues for accessing them (such as positions in a list or spatial locations)^{29,110}. Interference arises because the context representations of different items are similar to each other. For instance, the context for the first list item is similar to the context for the second item, so that when a person tries to retrieve the first item using the "position one" context as retrieval cue, the retrieved content is a blend of all list items, weighted by the similarity of their list position with the first position. The more the retrieved information is distorted in this way relative to the original information, the harder it is to recover the original from it. The interference hypothesis implies that the capacity of working memory depends on the confusability of the contents to be held and of the contexts to which they are bound, and on other variables determining the recoverability of distorted traces. Therefore, this hypothesis does not entail a natural unit of measurement for the capacity of working memory.

Oberauer and colleagues¹⁰² reviewed the evidence speaking to these competing hypotheses. They concluded that decay does not contribute to the capacity limit of working memory, and that although neither a limited resource nor interference are fully satisfactory explanations of the capacity limit on their own, a combination of both hypotheses – though not yet fleshed out as a theory – appears promising.

4.2. The relation of working memory and long-term memory

Working memory and LTM are related in two ways. First, knowledge in LTM facilitates maintenance and processing of information that corresponds to that knowledge in working memory. Second, new information that is encoded into working memory is also encoded – though perhaps only weakly – into episodic LTM, and these new representations in LTM can in turn assist the current task. We will review both aspects in turn.

4.2.1. Long-term knowledge assisting working memory

In his famous article on the capacity limits of the mind,⁶⁵ pointed out that we can briefly remember a list of words (e.g., "bat, ring, fan") much better than an equally long list of arbitrary letter strings (e.g., "bir, fong, ras"). Miller argued that known words are represented as a single unit – a chunk – in LTM and working memory, whereas arbitrary letter strings are lists of several units (i.e., the individual letters). More generally, when knowledge enables us to package information into larger units, keeping that information in working memory is easier. For instance, with increasing chess expertise, chess players excel more in reproducing the positions of pieces on a chess board after a brief glance, but only if the pieces are arranged in a way that could emerge from a chess game, because these arrangements contain many typical sub-configurations that chess experts are highly familiar with, so that they probably represent them as chunks^{111,112}.

One mechanism through which knowledge helps working memory is redintegration^{113,114}, the process by which the original stimulus is recovered from a distorted or corrupted memory trace. For instance, when the memory trace of "ring" is diminished to "r—g", the impoverished trace can be compared to all known words in the language, and "ring" is likely to be recovered as the best match. Obviously, this process works only if "ring" is known as a lexical unit. Immediate recall of word lists is influenced by a number of aspects of our lexical knowledge, including word frequency¹¹³, concreteness¹¹⁵, and the number of orthographic neighbors (i.e., words differing from the target word by only one letter)¹¹⁶.

These effects can be understood, at least partly, as reflecting the influence of knowledge on redintegration.

The beneficial effect of chunking on working memory, however, goes beyond the facilitation of redintegration. When people are asked to remember lists of letters or words in which some but not all of the list items can be integrated into larger chunks (e.g., the list FBIDKA, in which the first 3 letters form a chunk but the last 3 don't), their recall performance exceeds that of an unchunked control condition not only for the chunked items but also for the unchunked items on the list¹¹⁷. Because an effect of knowledge on redintegration could only help recovering the chunked, not the unchunked part of the list, it appears to be the case that chunking of part of the information in working memory frees capacity for other information.

4.2.2. Long-term memory used for short-term remembering

Performance on working memory tasks can also be influenced by new knowledge acquired during a working-memory testing session. On tests of immediate serial recall of short lists of items, for example, when the same list is repeated every third trial, performance improves rapidly for reproduction of the repeated list. This is observed even though subjects are not told to remember the lists for longer than a single trial⁶⁷. This so-called Hebb effect demonstrates that LTM acquires some information on every trial of a task designed to measure working memory, and that the resultant accumulation of knowledge across trials contributes to performance.

One inescapable implication of this finding is that at least some of the tasks commonly used for investigating working memory – and perhaps all these tasks – are not process pure: Performance on these tasks reflects a mixture of contributions from working memory and from rapidly acquired LTM. One strong interpretation that has been drawn from this is that the distinction between working memory and LTM is artificial, and that there may be only a single memory system responsible for maintaining information over any time scale between seconds and years^{118,119}. On this unitary-memory view, working memory is best described as the recruitment of general memory mechanisms for maintaining efficient access to very recently used information.

One strong argument for the unitary view is that it has been very difficult to demonstrate a convincing double dissociation between working memory and LTM. A classical neuropsychological dissociation relies on the observation that patients with damage to the hippocampus are severely impaired in acquiring new explicit LTM whereas their performance in tests of working memory is usually unimpaired (for a review see¹²⁰). There are, however, exceptions to this dissociation: Performance on some tasks presumed to test working memory has been found to be impaired in people with lesions to the hippocampus; this deficit appears to be specific to the ability to form and maintain bindings between objects and their spatial locations^{121,122}. Hence, one core function of working memory – maintenance of temporary bindings between contents and their contexts – appears to rely in part on the hippocampus, at least in some cases involving spatial context, rendering the neuropsychological dissociation between working memory and episodic LTM less clear-cut than would be desirable. A second limitation of the neuropsychological dissociation is that, whereas there are numerous reports of patients with selective deficits in tests of working memory for specific contents such as phonological information or spatial information (for a

brief review see ¹²³), no cases have been reported with a selective deficit in general working memory but not LTM. As long as selective working-memory deficits are limited to specific content domains, it is likely that they reflect an impairment of representations in that domain (e.g., a corruption of phonological codes in the mental lexicon) rather than of the mechanisms for holding representations available for guiding cognitive processes.

A second dissociation between episodic LTM and working memory that has been proposed is that episodic LTM is vulnerable to proactive interference whereas working memory is not ^{124,125}. This general claim cannot be upheld in light of evidence showing proactive interference in immediate tests of memory using very small set sizes – conditions that undoubtedly maximize the involvement of working memory ¹²⁶. A revised version of the original hypothesis, however, might still be viable: Whereas information in working memory is vulnerable to proactive interference from items that had themselves been held in working memory during previous trials from the same testing session, information in working memory may nonetheless be shielded against proactive interference from contents of LTM acquired prior to the testing session ⁶⁸. This hypothesis is consistent with our functional analysis of the requirements of an efficient working memory (*section 2.3*), but it has not yet been thoroughly tested.

4.2.3 Is working memory different from activated long-term memory?

As summarized in *section 3.2*, the idea of working memory functions arising from an activated state of LTM representations is central to state-dependent models. At the theoretical level, this remains a hotly debated proposition ^{79,127-129}. There is also a considerable amount of research from cognitive neuroscience that is relevant to the relation of working memory to LTM, and this will be considered in *section 5*. Before leaving this topic, we'll consider one promising way forward that is suggested by a computational model of episodic memory that incorporates aspects of both views.

In Farrell's Temporal Clustering and Sequencing model¹³⁰, the continuous stream of events that we experience is organized in memory into hierarchically embedded episodes. The model's architecture is as in Figure 2B, using a hierarchy of embedded event contexts. Events belonging together in an episode are tied together by being bound to a common context. For instance, in an experiment asking participants to remember lists of words for immediate recall, each trial would form one episode, so that all words in that trial's list are bound to the same list context. Within a list, subsets of 2-5 words are encoded as groups that form smaller episodes embedded in the list episode. Retrieval of an episode is usually a two-step process: The first step is to access the context of the to-be-retrieved episode; the second step is to use this context as a retrieval cue to the events bound to it. The last-encoded episode (for instance, the last group of a list of words) is assumed to have a special status in memory because its context is still active, so it does not have to be retrieved. We could think of the most recently experienced set of events as the contents of working memory. They are particularly well accessible because access to them does not require an error-prone retrieval of the relevant context. Because proactive interference arises mainly at retrieval of the episodic context, access to the contents of working memory is largely shielded from proactive interference (but see *section 5.4.2*). By the assumption that hippocampal damage primarily impairs context retrieval, the model can also explain why damage to the hippocampus tends to spare memory for the events in the most recent episode.

4.3. The relation of working memory to attention and cognitive control

Although most theorists assume that working memory is closely related to attention, this issue is complicated by the fact that many differ in how they conceptualize attention, and how they characterize its relation to working memory (for a review see ¹³³). For instance, in the multi-component model (*section 3.1*), the Central Executive component is a mechanism for what is often called "executive attention", that is, people's ability to control their own thoughts and actions to keep them aligned with their current goals. In Cowan's ¹³¹ embedded-process model (*section 3.2*), the focus of attention is characterized as a limited attentional resource that is needed for maintaining up to about four chunks in a highly accessible state. In neuroscience research, as we shall see in *section 5*, much of the emphasis is on the mechanisms, and effects, of sensory/perceptual attention.

One way of conceptualizing the relation of working memory and attention is to assume that attention is a limited mental resource that is responsible for the capacity limit of working memory. An alternative conceptualization describes attention not as a resource but as a collection of mechanisms for selectively prioritizing some information for processing. We next review theoretical ideas and evidence pertaining to these two perspectives.¹³²

4.3.1. Attention as a resource

The idea that a limited attentional resource is needed to maintain information in working memory has been fleshed out in three different ways: (1) A resource for short-term storage and processing, (2) a shared resource for perceptual attention and short-term maintenance, and (3) a resource for cognitive control.

4.3.1.1. An attentional resource for storage and processing.

The idea that short-term maintenance and processing of information must share a limited resource has a long history^{133,134}. Its most recent installment is the time-based resource-sharing (TBRS) theory¹³⁵. The TBRS theory starts from the assumption of a bottleneck for central cognitive processes, such as making a decision about how to respond to a stimulus, or retrieving information from LTM¹³⁶. This bottleneck is assumed to be required for refreshing representations in working memory that would otherwise decay. When additional cognitive processes are required during the retention interval of a working-memory task, these processes compete with refreshing for the central bottleneck. The TBRS theory points to this competition to explain why memory performance declines monotonically as the temporal density of concurrent processing demands is increased.

One problem for this theory, however, is that although it predicts competition between refreshing and a concurrent processing demand throughout the retention interval, this appears not to be the case: The effect of memory load on the speed of a concurrent processing task diminishes rapidly over the first 2-3 seconds of the retention interval¹³⁷⁻¹³⁹, and sometimes disappears completely after a few seconds¹⁴⁰⁻¹⁴².

4.3.1.2. Perceptual attention and working memory

Research on working memory for visual and spatial information has revealed a high degree of overlap between attention to perceived visual stimuli and maintenance of no-longer visible stimuli in working memory (we review evidence from neuroscience concerning this relationship in *section 5.2*). If perceptual attention is conceptualized as a resource, this overlap suggests that the same resource is also demanded by working memory. Support for

this hypothesis comes from studies showing that people's ability to simultaneously attend to multiple visual objects is limited in a way very similar to their ability to maintain multiple visual objects in working memory^{143,144}. A shared resource between perceptual attention and working memory would lead to substantial dual-task costs when a task demanding perceptual attention is inserted in the retention interval of a working-memory task. Evidence for this prediction is mixed: Some studies have found that a load on working memory impairs performance on a perceptual-attention task¹⁴⁵, whereas others have found little, if any dual-task cost¹⁴⁶. Dual-task costs of combining memory loads with perceptual-attention demands appear to be larger when there is representational overlap between contents of working memory and the stimuli for the perceptual-attention task. For instance, working memory for visual objects is impaired more by a concurrent task involving attentional selection of objects whereas memory for spatial locations is impaired more by a concurrent visual-search task^{147,148}. Similar patterns of dual-task interference in the absence of perceptual overlap have been interpreted as evidence for multiple encoding in working memory¹⁴⁹. Passive viewing of¹⁵⁰ or listening to¹⁵¹ nouns, and making syntactic judgments about written words¹⁵², all selectively disrupt delayed recognition of nonrepresentational shapes, suggesting that working memory for these visually presented stimuli engages linguistic and/or semantic codes, in addition to perceptual ones. Conversely, self-generated eye movements made in the dark, with no visible targets, selectively disrupts delayed recognition of locations^{150,151}, suggesting that working memory for locations engages covert oculomotor codes, in addition to perceptual ones. What remains to be determined conclusively, however, is the extent to which these patterns of content-specific interference may reflect interference between memory representations and distracting stimuli/actions, versus competition for a shared attentional resource.

4.3.1.3 *Controlled attention and working memory*

Some researchers have argued that the capacity of working memory is closely related to people's ability to control their cognitive processes, keeping their attention focused on what is relevant for their current goal and avoiding distraction^{153,154}. This idea is often expressed in terms of a shared resource for working memory maintenance and cognitive control¹⁵⁵⁻¹⁵⁷.

Evidence speaking to this assumption comes from correlational studies: Many studies testing large samples of young adults have found that performance on working-memory tasks correlates with indicators of cognitive control, such as the size of the Stroop effect or the flanker effect, the efficiency of stopping an action in the stop-signal task, or the ability to move the eyes away from a flashing stimulus in the anti-saccade task (for reviews see^{153,158}). One problem with this line of research, however, is that multiple indicators of cognitive control often don't correlate well with each other, implying that they may measure task-specific skills rather than a general ability to control one's thoughts and actions^{158,159}. A second source of evidence speaking to the hypothesis of a shared resource for working memory and cognitive control is dual-task studies combining a working-memory maintenance task (e.g., remembering a list of digits) with a demand on cognitive control (a Stroop or flanker task). The assumption of a shared resource entails the prediction that a higher memory load leads to impaired cognitive control (e.g., larger Stroop or flanker effects, or increased susceptibility to irrelevant distractor stimuli). However, the evidence on this prediction is inconsistent: Some studies have found the predicted impairment in cognitive control^{160,161}; others have found the opposite --less distraction under higher cognitive load^{162,163}; and yet others have found that memory load can both increase and decrease

indicators of cognitive control, depending on which kind of stimuli are used in the two tasks¹⁶⁴⁻¹⁶⁷.

4.3.2 Attention as a selection mechanism

A second perspective on the relation between working memory and attention starts from the definition of attention as a set of mechanisms and processes by which the cognitive system prioritizes some of the information available from perception and memory for processing. From this point of view attention is not a limited resource – rather, the limit on what we can attend to at any time arises from the function of attention: Selective prioritization necessarily implies exclusion of most available information; attending to many objects or events at the same time undercuts the purpose of selective attention.

Building on this definition of attention, we can characterize working memory as a form of attention: The contents of working memory are the representations that are currently most available for processing, and as such they are prioritized over all other representations. Perceptual attention plays a role in controlling which sensory information is gated into working memory. Analogously, we can think of retrieval of information from LTM into working memory as a form of selective attention to memory (see *section 2.3.1*).

Attentional selection of memory representations appears to occur over several levels of increased selectivity, so that the contents of working memory can be described as embedded sets of representations, as envisioned in state-based theories of working memory: Within a large set of representations currently activated in LTM, a subset of about 2 to 6 chunks is selected for being in a highly accessible state, referred to as the (broad) "focus of attention"¹³¹ or the "region of direct access"⁴⁷. Within that set there might be a further level of selection when an even smaller subset – often a single chunk – is selected for processing by a (narrow) "focus of attention"¹⁶⁸.

Information from perception and memory is selected for a purpose: The contents of working memory are selected either as the objects of processing (e.g., holding in mind an intermediate product while performing mental arithmetic), or as the information needed for controlling cognitive processes. We have already discussed (in *section 2.2.2*) one way in which working-memory contents control cognition: Working memory holds procedural representations – the currently relevant "Plan" or task set – that controls how the declarative contents of working memory are processed. In addition, the declarative contents of working memory also serve a role in controlling cognition. This role can be illustrated by research on visual search: Searching for an object in a cluttered scene requires holding a template of the search target in working memory. Once the template representation is in working memory, it guides perceptual attention automatically to objects in the scene that match the template. This "attentional capture" effect has been demonstrated in numerous experiments in which participants are asked to hold a simple visual object (e.g., a red disk) in working memory for a subsequent memory test. During the retention interval an unrelated visual-search task is carried out in which, on some trials, one of the distractors matches the object in working memory. This leads to slowed search, and an increase of eye fixations on the matching distractor, indicating that the distractor matching the current content of working memory attracts attention even when this is detrimental to efficient visual search^{169,170}.

5. Neural bases of working memory

Working memory has been a focus of intensive research by neuroscientists for decades. Here we will review current thinking about how working memory is accomplished by the brain. The mapping between the cognitive and theoretical constructs that we have reviewed up to this point, and the neural data that will feature in this section, is rarely one-to-one. Furthermore, because the neuroscientific study of working memory has largely been carried out in a distinct scientific framework, a conceptually coherent review of it cannot mirror the organization of the preceding sections. Therefore, to help relate the content from the preceding sections of this chapter with what's to follow, sub-section headings for *section 5* will be annotated with terms pointing to the relevant concepts highlighted in the excerpt from Miller et al. (1960)¹⁷ that opened this chapter, and from *section 2: Functional Requirements for Working Memory*, as described in *Table 1*.

Table 1. Annotation terms for concepts from Miller et al. (1960)¹⁷ and from *Section 2*

annotation	concept
from Miller et al. (1960) ¹⁷	
system/state	whether working memory is better understood from a memory-systems or a state-dependent perspective
plans/data	the distinction between the rules guiding behavior versus the storage of situation/trial-specific information, and interactivity between these two levels of representation (note that the plans vs. data distinction corresponds to the concept, from <i>section 2</i> , of procedural working memory vs. declarative working memory)
LTM	the role of LTM in working memory
from <i>Section 2</i>	
binding	The temporary binding between stimulus information and its context
attention	Selective attention to elements in working memory
stability/flexibility	The controlled encoding into working memory of only those elements in the perceived environment that are relevant for the current task, and updating via the selective removal of a subset of the contents of working memory, often entailing its replacement with new information

5.1. 20th century study of the working-memory functions of the PFC [*systems/state; plans/data; attention*]

Although neuroscientists have made remarkable progress in our understanding of the neural bases of working memory functions, there remains a noteworthy lack of consensus about some fundamental questions, such as the role of circuits in the PFC in the storage of information, and the importance of elevated, sustained neural activity for the storage of information. A brief historical review will be helpful for the interpretation of the current literature.

5.1.1. Lesion studies

It is surprisingly, and somewhat disconcertingly, common in the contemporary literature to find authors motivating or otherwise contextualizing current work by citing a

single seminal experiment from the 1930s – that of Jacobsen (1936)¹⁷¹ -- but then neglecting to reference any of the several ensuing studies that require a qualification, if not an outright revision, of Jacobsen’s original interpretation. For this reason, this subsection will go into a more granular level of detail than is characteristic of the rest of this chapter.

For his influential study, Jacobsen¹⁷¹ trained two nonhuman primates (NHPs) to perform a delayed-response task in a variant of the Wisconsin General Testing Apparatus (WGTA). After watching while one of two covered food wells was baited, the animal was made to wait for several seconds during which a lowered screen blocked it from seeing or reaching the wells. When the screen was raised, the animal was given one reach with which to displace the cover and, on a correct response, retrieve the food. Pre-lesion, the animals learned to perform the task almost perfectly. After recovery from bilateral surgical removal of the prefrontal cortex (PFC) anterior to the arcuate sulcus, the animals’ performance never deviated from chance.

Although Jacobsen¹⁷¹ concluded that the prefrontal cortex is responsible for “immediate memory,” the idea that the storage, per se, of to-be-remembered information had been disrupted by damage to the PFC was ruled out by a series of studies carried out over the next three decades. In one, Malmo (1942)¹⁷² replicated the basic procedure from Jacobsen’s experiment, but added the experimental factor of turning off the lights in the lab on one half of the trials. Remarkably, this simple manipulation had the effect of rescuing the performance of the PFC-lesioned NHPs, in that they performed correctly on roughly 85% of lights-off trials, despite still getting only 50% correct on lights-on trials. Malmo¹⁷² attributed his findings to an increased susceptibility to interference after bilateral PFC removal. Many studies that followed used tasks that demonstrated an important role for PFC in the control of behavior that is guided by the contents of working memory. A classic example is delayed alternation, a continuous task in which the animal is rewarded on each trial for selecting the one of two available stimuli (or locations, or actions) that it did *not* select on the previous trial. Although NHPs with PFC lesions are impaired on the standard version of the task -- when trials occur in an unbroken series with 5-sec intertrial intervals (ITI)^{173,174} -- it was later shown that this impairment was not due to an inability to remember information from the previous trial. To do this, Pribram and Tubbs (1967)¹⁷⁵ first replicated the impairment of PFC-lesioned NHPs with trials requiring alternating reaches to the right and to the left that were separated by 5-sec ITIs. They were then able to rescue performance to the level of control animals by simply increasing the ITI between each left-reach trial and the ensuing right-reach trial to 15-sec. Note that although lengthening a delay period would be expected, a priori, to increase demands on memory storage, the authors suggested that it improved performance of the PFC-lesioned animals by making it easier to parse their behavior into discriminable chunks.

In a different task, Pribram and colleagues (1964)¹⁷⁶ presented NHPs with an array of “junk” objects, and required them to first discover, by trial-and-error selection, which one covered a reward (“exploration strategy”), then to continue selecting this rewarded object until a criterion level of five consecutive correct choices was achieved (“exploitation strategy”), after which the experimenter baited another object (out of view of the animal), effectively requiring a switch back to the exploration strategy. At the beginning of each testing session, PFC-lesioned animals made more errors before achieving criterion with the first baited item, a pattern that could have been due either to forgetting what choices they had recently made, or by an impairment in shifting between explore and exploit strategies. Once

they achieved criterion, however, this ambiguity was resolved, because the PFC-lesioned animals then also perseverated on the exploit strategy longer than did temporal lobe-lesioned and control animals. That is, the impairment didn't result from impaired memory for choices, but, rather, from an impairment in using that information to successfully guide behavior. This pattern of impairment, qualitatively similar to that seen with delayed-alternation, would likely be interpreted in the current literature in terms of impaired processing of prediction errors (c.f., ^{177,178}).

Contemporaneous research being carried out in humans pointed to similar conclusions. Patients with PFC damage were reported to be unimpaired on forward digit span¹⁷⁹ and on delayed recognition of nonsense shapes drawn from an open set (i.e. no stimulus repeated during the testing session, ¹⁸⁰). The latter group of PFC patients was impaired on tests of delayed-response for other stimulus material (flicker frequency, color, tones, click frequency), but on each of these tests the stimuli were drawn from closed sets, meaning that stimuli repeated over the course of the testing session, thereby increasing the level of proactive interference relative to the test using an open set. On the Wisconsin Card Sorting Test, patients with lesions of the dorsolateral PFC were unimpaired relative to control subjects at learning the first sorting dimension –meaning that they could remember their previous incorrect choices and not repeat them -- but then made a disproportionate number of perseverative errors when the sorting dimension changed¹⁸⁰.

More recently, in a conceptual replication of Malmö (1942)¹⁷², humans with PFC-lesions were shown to be disproportionately impaired when distracting tone pips were played during the delay period of trials of delayed recognition of environmental sounds¹⁸¹. Importantly, a follow-up study in which the EEG was recorded during performance of the same task gave some insight into the PFC-dependent mechanisms underlying this impairment: The N1 component of the ERP to the sample stimulus was suppressed in PFC-lesioned patients; and middle-latency components of the auditory evoked potential (MAEP) to the distractors were larger for the PFC-lesioned patients¹⁸². The first result, mirroring what had previously been observed in a test of auditory selective attention in PFC-lesioned patients¹⁸³, was interpreted as underlying an impairment in the ability “to focus attention on task-relevant stimuli”¹⁸² (p. 173). The second result, because the MAEP reflects the initial cortical processing of the auditory signal, demonstrated an impairment of filtering distracting sensory information. Thus, the work of Malmö (1942)¹⁷² and of Chao and Knight (1998)¹⁸² suggest a role for the PFC in the function of input gating, as discussed in *section 2.3.1*. More recently, input gating has been modeled as a function supported by recurrent circuitry between PFC and the basal ganglia (e.g., ^{184,185}).

To summarize, the preponderance of lesions studies carried out in the 20th century has shown that the working memory functions of the PFC relate more closely to the control of working memory (in the case of input gating), and the control of behavior guided by the contents of working memory, than to memory storage per se. With regard to the organization of working memory, these studies suggest an anatomical distinction between the implementation of “Plans,” linked by this work to the PFC, versus the storage of trial-specific information. Although they do not speak directly to the question of whether working memory is better understood from a memory-systems or a state-dependent perspective, they do argue against models that posit a specialized role for the PFC in the storage of trial-specific information. This latter point is missed, of course, in reviews (e.g., ^{186,187,188}) that cite Jacobsen^{171,189}, but omit consideration of the work that followed.

5.1.2. Neurophysiology of memory-guided reaching

By the late 1960s, refinements in the ability to record neuronal activity from the brains of awake, behaving animals allowed scientists to begin designing studies intended to identify neural correlates of working memory processes. The majority of these studies targeted the PFC, because the integrity of this region had previously been demonstrated to be important for performance on these tasks. During delayed-response performance, Fuster and Alexander¹⁹⁰ found that many neurons in both PFC and the mediodorsal (MD) nucleus of the thalamus displayed elevated firing rates that spanned the duration of the delay period, which varied in length, unpredictably, within a range of 15 to 65 seconds. During delayed-alternation performance, Kubota and Niki¹⁹¹ observed two classes of task-related activity: neurons with elevated activity during the delay; and neurons that became active just prior to, and during, the response period. It is noteworthy that, in their contemporaneous interpretations of these findings, neither group interpreted these patterns of PFC activity as relating to the storage, per se, of information (see Postle¹⁹² for a more detailed treatment of these studies).

5.1.3. Neurophysiology of oculomotor delayed response

A series of studies carried out by Patricia Goldman-Rakic, Shintaro Funahashi, and their colleagues at Yale University during the 1980s and 1990s has had a remarkably enduring influence on thinking about the working-memory functions of the PFC. Goldman-Rakic worked within a memory-systems framework, assuming that circuits within the prefrontal cortex were crucial for the storage of information in working memory, as well as its manipulation¹⁹³. The procedure for their studies was adapted from methods for studying the visual system: first identify the tuning properties of a neuron, then observe how its activity may vary as a function of the manipulation of an experimental variable (in this case, impose a delay between sample and test). Results of these studies provided evidence for sustained delay-period activity in PFC neurons tuned for sample location¹⁹⁴⁻¹⁹⁷ or for sample identity¹⁹⁷, and were interpreted as evidence for a memory-storage function for the PFC. Evidence for a critical memory storage function was also seen in the fact that small, unilateral lesions of dorsolateral PFC produced impaired oculomotor delayed response -- but spared visually guided saccades -- to circumscribed locations in the contralateral visual field (an effect referred to as “mnemonic scotomas”¹⁹⁸). More specifically, Goldman-Rakic’s model posited a domain-specific organization of mnemonic function, with circuits in dorsolateral PFC responsible for the storage of location information, and ventrolateral PFC responsible for the storage of object information^{197,199,200}.

In the ensuing years, several studies have offered alternative interpretations to these findings, including: the seeming selectivity of PFC neurons may be a consequence of behavioral task and/or training^{201,202}; tasks that unconfound the focus of attention from the contents of working memory show PFC neurons to be more strongly related to the former²⁰³; and the “scotomas” produced by small unilateral PFC lesions may reflect greater susceptibility to proactive interference or to behavioral perseveration, rather than exaggerated forgetting²⁰⁴. Subsequently, there have been arguments raised that challenge these alternative interpretations (e.g.,^{205,206}). Nonetheless, it has been suggested that although Goldman-Rakic, Funahashi, and colleagues assumed that the sustained, stimulus-tuned activity they recorded from the PFC corresponded to the operation of the inner scribe and visual cache buffers from the multiple-component model of working memory, they may have, instead, been recording from neurons that contribute to the Central Executive¹⁹².

5.2. Circuit-level mechanisms of the control of visual working memory [*systems/state; plans/data; attention*]

Many neurally inspired state-dependent models assume that the retention of information in working memory relies on the mechanisms of selective perceptual attention (e.g., ^{207,208,209}). Thus, an important question for these models is whether the source(s) of the top-down control of spatial attention and of object-based attention play a similar role in visual working memory. For memory-systems models, in contrast, it is important to find evidence for specialized properties of neurons in higher-level regions of cortex that enable them to maintain information over a delay, and for differential patterns of connectivity in the cells and circuits responsible for the storage of information in working memory. There currently exists evidence consistent with both of these perspectives.

5.2.1. State-dependent models

5.2.1.1 *Spatial working memory*

Spatial selective attention is tightly linked to the circuitry that controls the direction of gaze. For example, after identifying the region of the visual field to which suprathreshold electrical microstimulation in the frontal eye field (FEF)^{210,211} or the superior colliculus^{212,213} will drive the eyes (i.e., a neuron's "motor field"), subthreshold microstimulation produces attention-like enhancement of detection of search targets at that location. Furthermore, this subthreshold microstimulation also produces attention-like enhancement of the visually driven response of V4 neurons with receptive fields overlapping the stimulated FEF motor field, enhancements that are greater for stimuli for which the V4 neuron is optimally tuned, and when a distractor is present elsewhere in the visual field²¹⁴.

If spatial working memory is believed to depend on sustained attention allocated to the to-be-remembered location(s) in space, one would expect, based on the findings summarized above, that spatial working memory also engages the circuitry involved in oculomotor control. There is, indeed, considerable evidence to support this proposition. In NHPs performing a task that required memory for a cued location, followed by a lever-release response (i.e., no eye movements throughout the trial), neurons with motor fields overlapping the cued location showed elevated activity throughout the delay period, and errors were associated with weakening of this activity. Furthermore, the remembered location could be decoded with remarkably high accuracy with multivariate pattern analysis (MVPA) of the full sample of recorded neurons²¹⁵. Even in the absence of an overt working memory task, during free viewing behavior, neurons in the FEF of NHPs encode information about recent saccade targets²¹⁶. Finally, pharmacological inactivation of the FEF with muscimol, a GABA_A agonist, devastates performance on a test of oculomotor delayed response, although it leaves object delayed match-to-sample performance relatively unaffected²¹⁷.

In humans, MVPA of fMRI activity from posterior superior frontal cortex (pSFC), a homologue of the NHP FEF, and from intraparietal sulcus (IPS, a region also implicated in spatial attention and oculomotor control) indicates that the neural encoding of an egocentrically defined location is highly similar whether subjects are engaged in planning a delayed saccade to a visible target at that location, covertly attending to this target in order to detect a change in its luminance, or preparing a delayed response to this location when it must be remembered across a delay. Specifically, a decoder trained to discriminate leftward vs. rightward oculomotor *intention* can decode the analogous information from the *attention* and *retention* tasks, and the same is true for the other two²¹⁸. (More on MVPA in *section 5.3.1.*) Furthermore this functionality is specific to pSFC, and specifically does not generalize to the more anterior regions of dorsolateral PFC emphasized in *section 5.1*, because damage to²¹⁹ and repetitive transcranial magnetic stimulation (rTMS) of²²⁰ dorsolateral PFC in

humans only disrupts spatial working memory performance when the pSFC is also affected by the intervention.

5.2.1.2 Object working memory

Although the neural bases of the source(s) of endogenous object- and feature-based attention aren't as well understood as are those of spatial attention, a region that is anatomically proximal to frontal oculomotor control circuits has been implicated in the control of object-based attention, and so may also be important for object working memory. This region, in posterior ventrolateral PFC, is known as the inferior frontal junction (IFJ, at the intersection of the inferior frontal and precentral sulci) in the human, and the ventral prearcuate area (VPA) in the NHP. In humans, Baldauf and Desimone (2014)²²¹ observed with magnetoencephalography (MEG) that alternating attention between superimposed streams of translucent images of faces and of houses produced the expected alternations of attention-related boosts of signal intensity in stimulus-related activity in posterior face- and house-sensitive regions, and these were tightly linked to alternations in the strength of coherence in the upper gamma band (roughly 60-100 Hz) between IFJ and these posterior regions. In the NHP, Bichot and colleagues (2015)²²² have demonstrated that, in a visual search task, neurons in VPA showed selectivity for the search target and showed feature-based attentional modulation earlier than did neurons in FEF. Furthermore, local inactivation of VPA neurons produced marked deficits in search performance, and abolished the feature-based attention modulation of FEF that was observed prior to the inactivation²²². Whether VPA might also play a role in object working memory is a question to which we will return.

5.2.2 Memory-systems models of circuit-level mechanisms

There is considerable evidence that the PFC does, indeed, have distinct properties that one would want in a specialized working-memory system. Circuits in the PFC have distinctive patterns of recurrent connectivity that support formation of dynamical attractors that can stably represent information across delay periods in the absence of sensory input (e.g., ^{223,224-226}). Furthermore, pyramidal neurons in the PFC have distinct morphological and physiological properties relative to early sensory areas, and PFC has different proportions of interneurons, all of which may give the PFC a unique ability to support sustained delay-period activity¹⁸⁷.

There have been several reports of stimulus-specific delay-period activity in PFC neurons that have been interpreted as evidence for a storage function, the majority, including those reviewed in *section 5.1.3*, have required a reach or an eye movement to a remembered location. One study that has reported evidence for stimulus-selective activity in PFC for a nonspatial visual features, the direction of global motion in a random-dot kinematogram (RDK; a.k.a. “dot motion”), has provided important data by recording simultaneously from three brain areas of the NHPs performing the task. In this experiment by Mendoza-Halliday, Torres, and Martinez-Trujillo²²⁷, an RDK was followed by the serial presentation of two probe stimuli, both presented at a different location on the screen, only one of which matched the sample. In visual area MT, which is specialized for the perceptual analysis of motion, robust sensory-related direction-selective activity dropped to baseline levels soon after sample offset. In two regions located downstream from MT – visual area MST and the PFC -- delay-period spiking patterns supported robust decoding of sample identity. To this, the authors applied a memory-systems interpretation, proposing “a functional boundary between early visual areas [including MT], which encode sensory inputs, and downstream association areas [including MST and PFC], which additionally encode mnemonic representations” (p. 1255)²²⁷. That is, the authors proposed that working memory representations are different from sensory representations. (Another important finding from this study, that we will revisit

further along in this chapter, is that although neurons in MT did not display elevated firing during the delay period, delay-period oscillations in the local-field potential (LFP) recorded from MT did carry information about the sample stimulus. Furthermore, a causal influence from PFC was seen in the form of elevated coherence between PFC spikes and the phase of the LFP in MT at lower frequencies, particularly in the beta band. The strength of PFC-MT spike-field coherence in the theta, alpha, and beta bands was markedly lower on error trials relative to correct trials.)

This idea of a segregation of mnemonic from sensory functions was reinforced in a follow-up study from Mendoza-Halliday and Martinez-Trujillo (2017) that reported evidence suggesting at least partial segregation of PFC neurons that preferentially represent visual features of stimuli while they are being perceived versus while they are being remembered²²⁸. Such evidence is important for memory-systems accounts of working memory, because it provides “a substrate for discriminating between perceptual and mnemonic representations of visual features” (p. 1)²²⁸.

5.3. The delay-period representation of information [*systems/state; attention; LTM*]

The neuroscience of working memory has been strongly influenced by Hebb’s (1949)²²⁹ articulation of a dual-code theory for the retention of information in the nervous system: (1) an initial activity-based code holds a record of the to-be-remembered information until (2) synaptic reorganization establishes the weight-based code that is the basis for LTM. Building on this idea, a guiding assumption in working memory research has been that storage depends on sustained, elevated activity in the circuits representing the to-be-remembered information. Research over the past decade has generated a large amount of data that has led to many refinements to, and in some cases reconsiderations of, this longstanding assumption.

5.3.1. EEG and fMRI correlates of delay-period activity

The decade of the 2000s witnessed two developments that have had profound influence on the cognitive neuroscience of working memory. Vogel and Machizawa (2004)²³⁰ recorded the EEG while subjects performed a variant of the change-detection task used to estimate visual working memory capacity (see *section 4.1.2.1*), in which a precue indicated which of two sample arrays presented simultaneously, one in each visual field, was relevant for that trial. Their finding was that the “contralateral delay activity” (CDA), derived by subtracting ipsilateral from contralateral signals from electrodes over posterior parietal and occipital regions, scaled monotonically with set size within the range from 1 item up to the individual’s capacity, and then flattened off such that it never exceeded the amplitude corresponding to that individual’s capacity. Thus, the CDA indexed the number of items that a subject held in working memory, rather than the number of items presented in the sample array. The CDA is covered in detail in the chapter from Foster, Awh, and Vogel¹⁰¹ in this volume. Of further interest here will be how the CDA indexes interactions between working memory and LTM, as well as an analogue of the CDA has been observed in fMRI studies that have identified a region of the intraparietal sulcus (IPS) for which fMRI signal intensity also scales with estimates of working memory capacity²³¹⁻²³³.

A second important development was the realization that because even single-subject neuroimaging datasets were high dimensional, neuroimaging data were amenable to “information-based” multivariate analysis methods adapted from machine learning: multivariate pattern analysis (MVPA; e.g. ²³⁴⁻²³⁶). The gist of MVPA is that, rather than aggregating across large numbers of voxels to extract a single value of the spatially averaged activity level in a region (as was done, for example, by ^{231,232,233}), one can train classifiers to assess whether the pattern of activity across all the voxels in a region is systematically

different for different stimuli (e.g., for different directions of motion²³⁷). Importantly, successful decoding of the contents of working memory from patterns of activity during a delay period does not require that the aggregate level of activity is different from baseline.

As we have already seen in the study of Mendoza-Halliday and colleagues²²⁷, successful decoding of delay-period activity provides evidence that signals from the area in question contain information about the stimulus being remembered. Early demonstrations of the insights to be gained from MVPA about working memory included a demonstration that delay-period signal during a working memory task (delayed paired-associate recognition) could be shown to reflect the temporary activation of information from episodic LTM²³⁸, and that early visual cortex, including V1, maintains active delay-period representations of sample information, despite the fact that aggregate levels of delay-period signal intensity may not differ from baseline^{239,240}.

5.3.2. Functional role(s) of delay-period stimulus representation in different brain areas.

5.3.2.1. Sensorimotor recruitment

Building on the findings from Harrison and Tong²³⁹ and Serences and colleagues²⁴⁰, several studies have demonstrated that the decoding of stimulus information from delay-period signals in early visual regions, including V1, is sensitive to manipulations of such factors as attention and load (e.g.,^{241,242,243}), and in a manner that covaries with behavioral indices of the precision of remembered information²⁴³⁻²⁴⁵. These findings have been interpreted as evidence for a “sensorimotor recruitment” mechanism supporting visual working memory, whereby the same systems that are involved in the sensory perception of information, as well as for the execution of actions tied to this information (e.g.,^{218,246}), contribute to the storage of this information^{247,248}. By this account, it is sensory representations in early visual areas that are the targets of top-down modulation from the frontal systems described in *section 5.2.1*.

Detailed consideration of the proposed mechanisms of sensorimotor recruitment indicates that, although the concept of “activated LTM” is contentious in the cognitive psychology literature (see *section 4.2.3*), it is accepted as a given in the neuroscience literature. This follows from the fact that visual object recognition and visual perception, more generally, depend on the interaction between the bottom-up processing of incoming sensory information and pre-existing representations of visual knowledge (i.e., LTM) – without this “activation of LTM” the perceiver would experience a visual agnosia. Importantly, there’s considerable evidence that real-time visual perception involves recurrent activity between multiple levels of visual processing (e.g.,²⁴⁹⁻²⁵²), and so sensorimotor recruitment can be understood as a prolongation of this interactive process. From this perspective, even patterns of sustained delay-period activity in early visual cortex^{227,237,243,244,253}, including V1^{239,240,254,255}, might be understood as a consequence (if not a demonstration) of activated LTM.

5.3.2.2. Memory-systems accounts of the storage of visual features.

Whereas the previous section emphasized early visual areas, many studies using multivariate analyses of fMRI data have also found evidence for the delay-period representation of low-level stimulus features (e.g., line or grating orientation, direction of motion, color) in parietal^{244,256-262} and frontal cortex^{257,260-262}, and interpretations of these findings vary. Xu appeals to the capacity-related activity in IPS introduced in *section 5.3.1*, and further argues that because stimulus decoding in occipital cortex is abolished by the concurrent presentation of distractors²⁵⁹, “early visual areas are unlikely to ... serve[...] as the primary storage site for [visual working memory storage] (p. 801)”²⁶³. Christophel and colleagues (2017)²⁵⁷ have proposed a “division of labor” (p. 494) whereby “sensory cortex

maintains a high-resolution representation of the currently attended memory item [i.e., the item immediately relevant for behavior], and parietal cortex has low-resolution representations of both attended and unattended items” (p. 496). (In particular, this parietal representation is proposed to reflect a “cortical specialization” (p. 494) for working memory storage that obviates the need to invoke the activity-silent mechanisms that will be considered in *subsection 5.3.3.2*. Some of the debate between proponents of sensorimotor recruitment and cortical specialization models can be found here^{187,263-266}.

5.3.2.3. Evidence for working memory-LTM interactions in delay-period signals.

As addressed in *section 4.2.3*, the question of whether the active maintenance of information in working memory involves more than the temporary activation of representations from LTM remains controversial. One study that addressed this question explicitly used the following reasoning: If holding information in working memory involves actively representing information in a buffer that is not engaged during retrieval-from-LTM tasks that don't make explicit demands on working memory, an MVPA classifier trained to decode information from such a retrieval-from-LTM task should fail to decode this information during a working memory task. To test this hypothesis, Lewis-Peacock and Postle (2008)²³⁸ trained classifiers to discriminate the categories *celebrities*, *famous locations*, and *common objects* from fMRI data acquired during an initial scanning session, while subjects viewed individually presented images of exemplars from the three categories and made Likert-scale judgments about them: “How much do you like this person?; “How much would you like to visit this location?; and “How often do you encounter this object in your daily life?”. This task required retrieval of information from LTM, but made no overt demands on working memory. Next, outside the scanner, subjects learned to associate arbitrarily selected pairs of exemplar stimuli. Finally, during a second scanning session, subjects performed a delayed paired associate-recognition task in which one of the two items from each pair was presented as the sample stimulus, and the probe either was or was not the sample's paired associate. Results indicated that the classifiers trained on data from the retrieval-from-LTM task were able to detect the active representation of the category of the item paired with the sample item throughout the delay period. This suggested that, upon seeing the sample stimulus, subjects prospectively activated a representation of the sample's associate and held it in working memory in anticipation of comparing it with the probe. Additionally, successful classification of this activity meant that performing the delayed paired-associate-recognition task must have generated the same patterns of neural activity as had performing the retrieval-from-LTM task²³⁸.

A second set of studies has made clever use of the CDA to generate neural evidence for an analogue of the Hebb effect (*section 4.2.2*) across repeated trials of a visual search task. The basic logic of these experiments was to present sets of one or two search targets in each visual field, the relevant item(s) indicated by color, and to separate the offset of the target(s) from the onset of the search array by a 900 msec delay period. This allowed for assessment of the delay-period retention of the lateralized target item(s), which served as one or two search templates. Analyses of an initial experiment that varied which visual field was cued and the number of search targets confirmed that this procedure produced a CDA during the delay period. In a second experiment the authors presented one search target in each visual field, and varied the number of consecutive trials on which the same stimulus was cued (3, 5, or 7 consecutive trials). The results indicated that the amplitude of the CDA declined across trials in a pattern that could be fit by a power-law function. The authors interpreted this result as evidence for a “handoff of the attentional template from visual working memory to long-term memory as subjects searched for the same target object across runs of trials” (p. 9320)²⁶⁷. A follow-up study replicated the effect that consecutive-trial repetitions of the

search target were associated with a decline in the CDA, and also showed that a different ERP component, the P170, increased with each repetition²⁶⁸. Because the P170 indexes the accumulation of information that supports recognition from LTM, this provided further evidence for the engagement of working memory and LTM on the same task.

An additional source of evidence for an influence of LTM on working memory has come from the application of ideas from dynamical systems theory. In tests of delayed recall of color from arrays of one versus three (for humans) or two (for NHPs) colored squares (modeled on the tasks discussed in *section 4.1.2*), recall responses have been observed to be markedly biased away from some colors and toward others. Attractor dynamics accounted for the frequency, bias, and precision of these responses, with the greater error on high-load trials shown to reflect both a drift of remembered stimulus representations toward stable attractor states and a greater influence of random diffusion (i.e., noise). The authors framed this as evidence for an error-correcting mechanism, whereby increased internal noise (manipulated here by varying load) is counteracted by drift toward stable long-term representations of color space²⁶⁹. (From a Bayesian perspective, one could construe this as drawing on prior knowledge to counteract uncertainty about the recently presented stimuli.) Applying this model to behavioral data acquired during an fMRI study of delayed recall of one versus three line orientations²⁷⁰ has revealed that drift and diffusion parameters from the discrete attractor model relate closely to load-related changes in IPS²⁷¹, thereby suggesting an alternative explanation for effects previously attributed to working memory storage²³¹⁻²³³.

5.3.3 Alternatives to sustained, elevated delay-period activity

In recent years there has been growing interest in the possibility that mechanisms other than sustained, elevated activity may underlie the retention of information in working memory.

5.3.3.1 Transient attractor states underlying “gamma bursts”

Computational models from Lansner and colleagues use short-term Hebbian plasticity, driven by afferent signals encoding sample stimuli, to create transient networks of PFC pyramidal neurons that encode stimulus information. Interactions between recurrently connected pyramidal neurons and inhibitory basket cells produce a regime whereby an “activity silent” state is imposed by default network oscillations in the beta band (roughly 15-35 Hz), with stochastically occurring (i.e., not periodic) brief narrow-band bouts of oscillation in the gamma band (roughly 40-100 Hz), prompting bursts of activity in these stimulus-encoding networks²⁷²⁻²⁷⁴. This framework was applied to a data set recorded from the PFC of two NHPs performing change detection for sequences of two or three colored squares. Whereas trial averaging of the LFP data yielded patterns of prominent, sustained power in the beta band that spanned the delay period, single-trial analyses yielded evidence consistent with the models from the Lansner group: stimulus-related spiking occurred only sporadically, and was tightly coupled with brief bursts of narrow-band oscillations in the gamma band²⁷⁵. The conclusion of the authors, that working memory is supported by discrete bouts of firing and that evidence for sustained delay-period activity may be an artifact of trial averaging, has prompted considerable interest and debate²⁷⁶⁻²⁷⁸.

5.3.3.2 Neurally active and neurally silent representations in working memory

Independent of the computational framework described in the previous subsection, the idea of short-term synaptic plasticity as a basis for the short-term retention of information has been proposed in several theoretical and computational contexts (e.g.,^{97,98,279}), and as a possible explanation for patterns of activity observed in temporal cortex of NHPs performing

working memory tasks^{280,281}. In human working memory, an “activity-silent” mechanism (e.g., ^{282,283,284}) has been raised as a possible explanation for the fact that MVPA evidence for a stimulus in working memory drops to baseline when a retro-cue informs the subject that that stimulus will not be relevant for the impending memory test. Importantly, on trials when a second, subsequent retro-cue indicates that this stimulus is again relevant for behavior, MVPA evidence for it returns²⁸⁵⁻²⁸⁸. Furthermore, when a pulse of transcranial magnetic stimulation (TMS) is delivered to sensory or parietal cortex, it produces a transient reactivation of MVPA evidence for the “silent” item in the concurrently recorded EEG, suggesting that, despite the absence of an active trace, the representation of this stimulus remains in an accessible state²⁸⁷.

In the wake of these retro-cuing studies, neural network models of several working memory tasks, including the retro-cuing task just described, have suggested that synaptic weight-based stimulus representation may be the default mechanism for delay-period stimulus representation, whereas its supplementation by activity-based stimulus representation is more context dependent^{289,290}. In simulations of the retro-cuing task, active representation drops off in these models when an item is deprioritized, mirroring the effects from fMRI and EEG studies²⁸⁵⁻²⁸⁸. In one, the delivery of a nonspecific pulse of activation to the network, simulating the pulse of TMS from the TMS-EEG study²⁸⁷, produces a brief return of decodability of the uncued item²⁹⁰. Simulations with recurrent neural networks have also demonstrated that the level and complexity of delay-period activity depends on the requirements of the task, even across tasks that begin with the same to-be-remembered information. For example, Orhan and Ma (2019)²⁹¹ demonstrated that for tasks with fixed delay durations, with recognition probes (i.e., where the motor response can't be known prior to the probe onset), and for which the network had prior training on other tasks, individual units exhibited only relatively brief bouts of elevated activity, whereas for recall tasks (i.e., where response is indicated by the sample) and tasks with variable delay durations, individual units exhibited more persistent activity. In the simulations of Masse and colleagues (2019)²⁸⁹, short-term synaptic plasticity was sufficient to support performance on tasks requiring recognition of the sample as it had been presented, but on tasks requiring transformation of this information, guidance of performance with a post-sample instruction cue, or discriminating distractor repetition from sample repetition, higher levels of persistent activity were observed in the network. These computational findings may help explain qualitatively similar observations from neurophysiological studies, which have documented that the stability of single unit activity varies with such factors as the predictability of the duration of the delay period, and the requirement to transform the sample to generate the appropriate response²⁹². Masse and colleagues (2019)²⁸⁹ apply this reasoning as a possible explanation for why delayed recognition of a direction of motion involved elevated delay-period activity in MST in the experiment of Mendoza-Halliday et al. (2014)²²⁷, which presented the test stimuli in a location different from that of the sample, but not in the downstream lateral intraparietal area (LIP) in studies using a similar task^{293,294}, but in which the test stimulus appeared in the same location as had the sample. We note that the same reasoning might also explain the discrepancy between the finding of sustained elevated activity in PFC in the study of Mendoza-Halliday et al. (2014)²²⁷, but the failure to find such activity in another study of delayed recognition of the direction of motion²⁹⁵.

5.4. Neural bases of the control of working memory [*plans/data; binding; stability/flexibility*]

As highlighted in *Section 5.1.1*, the deficits in working memory performance that result from damage to the PFC are best understood as disruptions in the ability to control behavior with information held in working memory, rather than as disruptions of the ability

to store information in an accessible state for brief periods of time. Here, we will relate the functioning of frontal and parietal cortex to factors, emphasized in *Section 2* and elsewhere, that are critical for successful performance on tests of working memory: the binding of content to context, and the ability to control interference.

5.4.1. Context binding

Recently, Gosseries, Yu, and colleagues (2018)²⁴⁴ explored the idea that the representations of the location of objects that are found in LIP and FEF of the NHP and in IPS and pSFC of the human²⁹⁶ might play an important role for binding the contents of visual working memory (i.e., visual objects) to contexts (i.e., their representation in a frontoparietal priority map). Furthermore, they reasoned that because delay-period stimulus representation is more prominently and reliably represented in occipitotemporal regions than in frontoparietal regions (e.g.,^{237,243}), delay-period activity in the latter might reflect, at least in part, the maintenance of content-to-context bindings. To test these ideas, they scanned subjects performing delayed recall for one direction of motion (“1M”), three directions of motion (“3M”), or one direction of motion and two colors (“1M2C”). Samples in the 3-item trials were presented serially, with a digit accompanying the response dial indicating whether the first, second, or third item to be presented was to be recalled. The 1M2C condition was particularly diagnostic because it required maintenance of the same number of items as the 3M condition but put a lower demand on content-context bindings: Memory for the ordinal position of the motion sample was not required on 1M2C trials, because it was the only item that a motion probe could interrogate; on 3M trials, in contrast, the correct binding of each sample to its context was necessary for recalling the correct item.

Results indicated that delay-period activity in IPS was markedly higher for 3M than 1M trials, but equivalent for 1M2C and 1M trials. The fact that 1M2C trials matched the 3M trials for the number of items to be held in working memory suggests that the load effect that has frequently been observed in IPS (e.g.,^{231,232,233,243,271}) may reflect, at least in part, varying demands on context binding, rather than varying demands on stimulus representation per se. Whereas the critical context in the study of Gosseries, Yu, and colleagues (2018)²⁴⁴ was ordinal position of an item in the presentation sequence, the pattern of effects that they reported has been replicated in a study in which the three sample items (either three oriented bars or one bar, one color, and one luminance patch) were presented simultaneously and in different locations on the screen²⁷⁰. Therefore, an important question for future research is how the brain differentially represents spatial context versus ordinal-position/temporal context.

5.4.2. Proactive interference

The control of proactive interference in verbal working memory has been studied extensively with Monsell’s (1978)²⁹⁷ “recent negatives” variant of delayed recognition, in which a small number of consonant letters (e.g., four) is presented as samples, and on a subset of trials the nonmatching (a.k.a. “negative”) probe is drawn from the samples presented on the previous trials. Subjects are slower to reject such recent-negative probes than non-recent negative probes, and they false alarm to them at an elevated rate. Studies with PET imaging²⁹⁸ and with brain-damaged patients²⁹⁹ have implicated left inferior PFC as playing an important role in the control of proactive interference in working memory for verbal material, and studies with fMRI³⁰⁰ and with repetitive (r)TMS³⁰¹ indicate that this control is applied at the time of the recognition decision.

Additional study suggests that the control operation carried out by the left inferior PFC is one of evaluating item context. A behavioral experiment using a response-deadline procedure indicated that the false-alarm rate to recent-negative probes is highest with very

short response deadlines, suggesting that the fast familiarity signal generated by the visual processing of this probe triggers a “match” response before the slower recollection signal containing contextual information can influence the decision. Consistent with a role for left inferior PFC in comparing familiarity versus recollection signals, rTMS delivered to this area “early” (from 0-250 msec after probe onset), but not “late” (from 500-750 msec after probe onset), produced an elevated false-alarm rate to recent-negative probes³⁰².

5.5. Cognitive Neuropsychology of Working Memory [*system/state*]

The goal of the research that we have summarized up until this point has been to investigate neural mechanisms underlying different aspects of working memory. A different approach that has contributed importantly to working memory theory development is the use of case studies from the neurology clinic, and sometimes larger samples of patients, to identify patterns of behavioral deficits that provide evidence for particular aspects of models of the cognitive architecture of working memory. Indeed, a report of patients with a selective impairment of auditory-verbal short-term memory (Warrington and Shallice, 1969)³⁰³ was influential in the development of modality-specific buffers in the multiple-component model⁷⁵. A thorough review of this literature is beyond the scope of this chapter, but its history is succinctly captured by Papagno and Shallice (2019)¹²³, and the special issue that this paper introduces provides a thorough overview of the current state of this literature, with arguments in favor of both memory-systems³⁰⁴⁻³⁰⁹ and state-dependent^{307,309,310} models of working memory.

6. Conclusion

Working memory is a construct that is of central importance for understanding many aspects of high-level cognition. Its study has generated influential discoveries and novel ideas in many branches of cognitive psychology and of neuroscience. In this chapter we have considered core functional and neural properties of working memory, providing an overview of the field by summarizing how currently influential perspectives address each of these properties. This exercise has highlighted the fact that increasing the crosstalk between cognitive/theoretical and neuroscientific approaches to the study of working memory offers an important way to make further progress. We conclude this chapter by considering some possibilities, framed by the concepts summarized in *Table 1*.

System/state; attention: Delay-period spiking in PFC coordinating the synchrony of LFPs in MT, as described by Mendoza-Halliday et al. (2014)²²⁷ (*section 5.2.2*), is an example of top-down control of activity in sensory cortex, the mechanism that is also widely assumed to be the basis of visual selective attention (e.g., ^{221,311,312,313}; *section 5.2.1.2*). Therefore, an experiment directly assessing evidence for overlap between the region of caudal PFC implicated in feature-based working memory²²⁷ and VPA, the region identified as a source of object-based attention²²² (*section 5.2.1.2*), could adjudicate between memory-systems and state-dependent accounts of activity in the PFC.

Plans/data: Are stimulus-specific patterns of activity observed in non-sensory cortex (such as the PFC; see *section 5.3.2*) best construed as supporting a storage function or a control? (For one study that has addressed this question, albeit in the context of visual perception rather than working memory, see³¹⁴.)

LTM: Does the neural representation of stimulus information being manipulated in working memory differ meaningfully from the neural representation of that same stimulus when it has been retrieved from LTM, but not in the context of a working memory task?

Binding: Does the binding of content to context that is fundamental to working memory (*section 2.1*) engage mechanisms that are different from the operation of the priority map hypothesized to govern visually guided behavior (*section 5.4.1*)?

Attention: Does the CDA that is measured during working memory tasks (*section 5.3.1*) reflect the operation of a mechanism that is qualitatively different from object-based attention?

Attention; Stability/flexibility: Is information that is in working memory, but not relevant for immediate behavior, represented differently from prioritized information?

Stability/flexibility: Advances in our understanding of how the principles of reinforcement learning are implemented in corticostriatal circuits (e.g.^{20,177,185,315,316}) can provide important constraints on our understanding of the mechanisms underlying prioritization and removal operations in working memory (e.g.^{285,317,318}).

Consciousness: Although the presence of activity in a particular brain area isn't specific to consciousness, it is reasonable to posit that activity is a necessary condition for consciousness to exist. Can further understanding of the distinction between activity-silent vs. active states of representation in working memory (*section 5.3.3.2*) provide useful leverage for understanding the neural bases of consciousness?

References

1. Shipstead Z, Redick TS, Hicks KL, Engle RW. The scope and control of attention as separate aspects of working memory. *Memory*. 2012;20:608-628.
2. Shipstead Z, Engle RW. Interference within the focus of attention: working memory tasks reflect more than temporary maintenance. *Journal of Experimental Psychology: Learning, Memory, & Cognition*. 2013;39:277-289.
3. Fukuda K, Vogel EK, Mayr U, Awh E. Quantity, not quality: The relationship between fluid intelligence and working memory capacity. *Psychonomic Bulletin & Review*. 2010;17:673-679.
4. Cowan N, Fristoe NM, Elliott EM, Brunner RP. Scope of attention, control of attention, and intelligence in children and adults. *Memory & Cognition*. 2006;34:1754-1768.
5. Ree MJ, Carretta TR. g2K. *Human Performance*. 2002;15:3-23.
6. Daneman M, Carpenter PA. Individual differences in working memory and reading. *Journal of Verbal Learning and Verbal Behavior*. 1980;19:450-466.
7. Gathercole SE, Pickering SJ. Working memory deficits in children with low achievements in the national curriculum at 7 years of age. *British Journal of Educational Psychology*. 2000;70:177-194.
8. Baddeley AD. *Working Memory*. London: Oxford University Press; 1986.
9. Jonides J. Working memory and thinking. In: Smith EE, Osherson DN, eds. *An Invitation to Cognitive Science*. Vol 3. Cambridge, MA: MIT Press; 1995:215-265.
10. Engle RW, Kane MJ, Tuholski SW. Individual differences in working memory capacity and what they tell us about controlled attention, general fluid intelligence, and functions of the prefrontal cortex. In: Miyake A, Shah P, eds. *Models of Working Memory*. Cambridge, U.K.: Cambridge University Press; 1999:102-134.
11. Devinsky O, D'Esposito M. *Neurology of Cognitive and Behavioral Disorders*. New York: Oxford University Press; 2003.
12. Gold JM, Barch DM, Feuerstahler LM, et al. Working memory impairment across psychotic disorders. *Schizophrenia Bulletin*. 2019;45:804-812.
13. Ebbinghaus H. *On Memory*. New York: English translation, Teachers' College (1913); 1885.
14. James W. *The principles of psychology*. New York: Henry Holt and Co.; 1890.
15. Feigenbaum EA, Simon HA. Performance of a reading task by an elementary perceiving and memorizing program. *The RAND Corporation Paper*. 1961;P-2358.
16. Newell A, Shaw JC, Simon HA. Elements of a theory of human problem solving. *Psychological Review*. 1958;65:151-166.
17. Miller GA, Galanter E, Pribram KH. *Plans and the Structure of Behavior*. New York: Henry Holt and Company; 1960.
18. Miller GA. What is information measurement? *American Psychologist*. 1953;8:3-11.
19. Oberauer K. Design for a working memory. *Psychology of Learning and Motivation*. 2009;51:45-100.
20. O'Reilly RC, Frank MJ. Making working memory work: A computational model of learning in the prefrontal cortex and the basal ganglia. *Neural Computation*. 2005;18:283-328.
21. Awh E, Vogel EK. The bouncer in the brain. *Nature Neuroscience*. 2008;11:5-6.
22. Halford GS, Wilson WH, Phillips S. Processing capacity defined by relational complexity: Implications for comparative, developmental, and cognitive psychology. *Behavioral and Brain Sciences*. 1998;21:803-864.
23. Shastri L, Ajjanagadde V. From simple associations to systematic reasoning: A connectionist representation of rules, variables and dynamic bindings using temporal synchrony. *Behavioral and Brain Sciences*. 1993;16:417-494.
24. Burgess N, Hitch GJ. A revised model of short-term memory and long-term learning of verbal sequences. *Journal of Memory and Language*. 2006;55:627-652.

25. Burgess N, Hitch GJ. Memory for serial order: A network model of the phonological loop and its timing. *Psychological Review*. 1999;106:551-581.
26. Brown GDA, Preece T, Hulme C. Oscillator-based memory for serial order. *Psychological Review*. 2000;107:127-181.
27. Lewandowsky S, Farrell S. Short-term memory: new data and a model. In: Ross BH, ed. *The Psychology of Learning and Motivation*. Vol 49. London, UK: Elsevier; 2008:1-48.
28. Oberauer K, Farrell S, Jarrold C, Pasiiecznik K, Greaves M. Interference between maintenance and processing in working memory: The effect of item-distractor similarity in complex span *Journal of Experimental Psychology: Learning, Memory, and Cognition*. 2012;38:665-685.
29. Oberauer K, Lin H-Y. An interference model of visual working memory. *Psychological Review*. 2017;124:21-59.
30. Schneegans S, Bays PM. Neural architecture for feature binding in visual working memory. *The Journal of Neuroscience*. 2017;37:3913-3925.
31. Henson RNA. *Short-term memory for serial order*, University of Cambridge; 1996.
32. Bays PM, Catalao RF, Husain M. The precision of visual working memory is set by allocation of a shared resource. *Journal of Vision*. 2009;9.
33. Pertzov Y, Dong MY, Peich M-C, Husain M. Forgetting what was where: The fragility of object-location binding. *PLoS ONE*. 2013;7. doi:10.1371/journal.pone.0048214.
34. Bays PM. Evaluating and excluding swap errors in analogue tests of working memory. *Scientific Reports*. 2016;6. doi:10.1038/srep19203.
35. Rerko L, Oberauer K, Lin H-Y. Spatially imprecise representations in working memory. *Quarterly Journal of Experimental Psychology*. 2014;67:3-15.
36. McElree B. Working memory and focal attention. *Journal of Experimental Psychology: Learning, Memory, & Cognition*. 2001;27:817-835.
37. McElree B, Doshier BA. Serial position and set size in short-term memory: The time course of recognition. *Journal of Experimental Psychology: General*. 1989;118:346-373.
38. Vergauwe E, Langerock N. Attentional refreshing of information in working memory: Increased immediate accessibility of just-refreshed representations. *Journal of Memory and Language*. 2017;96:23-35.
39. Vergauwe E, Hardman KO, Rouder JN, Roemer E, McAllaster S, Cowan N. Searching for serial refreshing in working memory: Using response times to track the content of the focus of attention over time. *Psychonomic Bulletin & Review*. 2016.
40. Oberauer K. Selective attention to elements in working memory. *Experimental Psychology*. 2003;50(4):257-269.
41. Garavan H. Serial attention within working memory. *Memory & Cognition*. 1998;26:263-276.
42. Donkin C, Nosofsky RM. A power-law model of psychological memory strength in short- and long-term recognition. *Psychological Science*. 2012;23:625-634.
43. Griffin IC, Nobre AC. Orienting attention to locations in internal representations. *J Cog Neuroscience*. 2003;15:1176-1194.
44. Landman R, Spekreijse H, Lamme VAF. Large capacity storage of integrated objects before change blindness. *Vision Research*. 2003;43(149-164).
45. Souza AS, Oberauer K. In search of the focus of attention in working memory: 13 years of the retro-cue effect. *Attention, Perception, and Psychophysics*. 2016;78:1839–1860.
46. Niklaus M, Singmann H, Oberauer K. Two distinct mechanisms of selection in working memory: Additive last-item and retro-cue benefits. *Cognition*. 2019;183:282-302.
47. Oberauer K. Design for a working memory. *Psychology of Learning and Motivation: Advances in Research and Theory*. 2009;51:45-100.
48. Anderson JR, Lebiere C. *The atomic components of thought*. Mahwah, N. J.: Erlbaum; 1998.
49. Koch I, Gade M, Schuch S, Philipp AM. The role of inhibition in task switching: A review. *Psychonomic Bulletin & Review*. 2010;17:1-14.

50. Kiesel A, Steinhauser M, Wendt M, et al. Control and interference in task switching - a review. *Psychological Bulletin*. 2010;136:849-874.
51. Vandierendonck A, Liefoghe B, Verbruggen F. Task switching: Interplay of reconfiguration and interference control. *Psychological Bulletin*. 2010;136:601-626.
52. Mayr U, Kliegl R. Task-set switching and long-term memory retrieval. *Journal of Experimental Psychology: Learning, Memory, and Cognition*. 2000;26:1124-1140.
53. Meiran N, Liefoghe B, De Houwer J. Powerful instructions: Automaticity without practice. *Current Directions in Psychological Science*. 2017;26:509-514.
54. Baddeley AD. Short-term phonological memory and long-term learning: A single case study. *European Journal of Cognitive Psychology*. 1993;5:129-148.
55. Colle HA, Welsh A. Acoustic masking in primary memory. *Journal of Verbal Learning and Verbal Behavior*. 1976;15(1):17-31.
56. Schlittmeier SJ, Weißgerber T, Kerber S, Fastl H, Hellbrück J. Algorithmic modeling of the irrelevant sound effect (ISE) by the hearing sensation fluctuation strength. *Attention, Perception & Psychophysics*. 2012;74:194-203.
57. Jones DM, Macken WJ. Phonological similarity in the irrelevant speech effect: Within- or between-stream similarity? *Journal of Experimental Psychology: Learning, Memory, and Cognition*. 1995;21(1):103-115.
58. Hughes RW, Vachon F, Jones DM. Disruption of short-term memory by changing and deviant sounds: support for a duplex-mechanism account of auditory distraction. *Journal of Experimental Psychology: Learning, Memory, and Cognition*. 2007;33:1050-1061.
59. Jones DM, Macken WJ. Irrelevant tones produce an irrelevant speech effect: implications for phonological coding in working memory. *Journal of Experimental Psychology: Learning, Memory & Cognition*. 1993;19:369-381.
60. Page MPA, Norris DG. The irrelevant sound effect: What needs modelling, and a tentative model. *Quarterly Journal of Experimental Psychology*. 2003;56A:1289-1300.
61. Bell R, Röer JP, Lang A-G, Buchner A. Distraction by steady-state sounds: Evidence for a graded attentional model of auditory distraction. *JEP:HPP*. 2019;45(4):500-512.
62. Unsworth N, Robison MK. The influence of lapses of attention on working memory capacity. *Memory & Cognition*. 2016;44:188-196.
63. Luria R, Balaban H, Awh E, Vogel EK. The contralateral delay activity as a neural measure of visual working memory. *Neuroscience and Biobehavioral Reviews*. 2016;62:100-108.
64. Oberauer K, Lewandowsky S. Simple measurement models for complex working memory tasks. *Psychological Review*. in press.
65. Miller GA. The magical number seven, plus or minus two: Some limits on our capacity for processing information. *Psychological Review*. 1956;63:81-97.
66. Hulme C, Maughan S, Brown GDA. Memory for familiar and unfamiliar words: Evidence for a long-term memory contribution to short-term memory span. *Journal of Memory and Language*. 1991;30:685-701.
67. Hebb DO. Distinctive features of learning in the higher animal. In: Delafresnaye JF, ed. *Brain mechanisms and learning*. Oxford: Blackwell; 1961:37-46.
68. Oberauer K, Awh E, Sutterer DW. The role of long-term memory in a test of visual working memory: proactive facilitation but no proactive interference. *Journal of Experimental Psychology: Learning, Memory and Cognition*. 2017;43:1-22.
69. Kessler Y, Meiran N. All updateable objects in working memory are updated whenever any of them are modified: Evidence from the memory updating paradigm. *Journal of Experimental Psychology: Learning, Memory & Cognition*. 2006;32:570-585.
70. Kessler Y, Oberauer K. Forward scanning in verbal working memory updating. *Psychonomic Bulletin & Review*. 2015;22:1770-1776.

71. Kessler Y, Oberauer K. Working memory updating latency reflects the cost of switching between maintenance and updating modes of operation. *Journal of Experimental Psychology: Learning, Memory, and Cognition*. 2014;40:738-754.
72. Rac-Lubashevsky R, Kessler Y. Dissociating working memory updating and automatic updating: The reference-back paradigm. *Journal of Experimental Psychology: Learning, Memory and Cognition*. 2016;42:951-969.
73. Miyake A, Shah P, eds. *Models of Working Memory*. Cambridge, U.K.: Cambridge University Press; 1999.
74. Logie R, Camos V, Cowan N, eds. *Working Memory: State of the Science*. Oxford University Press; in preparation.
75. Baddeley AD, Hitch GJ. Working Memory. In: Bower GH, ed. *The Psychology of Learning and Motivation*. Vol 8. New York: Academic Press; 1974:47-89.
76. Baddeley AD. The episodic buffer: a new component of working memory? *Trends in Cognitive Sciences*. 2000;4:417-423.
77. Baddeley AD. Fractionating the Central Executive. In: Stuss DT, Knight RT, eds. *Principles of Frontal Lobe Function*. Oxford: Oxford University Press; 2002:246-260.
78. Baddeley AD, Hitch G. Working memory: past, present ... and future? In: Osaka N, Logie RH, D'Esposito M, eds. *The Cognitive Neuroscience of Working Memory*. Oxford, U.K.: Oxford University Press; 2007:1-20.
79. Baddeley AD, Hitch GJ. The phonological loop as a buffer store: An update. *Cortex*. 2019;112:91-106.
80. Murray A, Jones DM. Articulatory complexity at item boundaries in serial recall: the case of Welsh and English digit span. *Journal of Experimental Psychology: Learning, Memory, & Cognition*. 2002;28:594-598.
81. Cowan N. The magical number 4 in short-term memory: A reconsideration of mental storage capacity. *Behavioral and Brain Sciences*. 2001;24:87-185.
82. Ungerleider LG, Mishkin M. Two cortical visual systems. In: Ingle DJ, Goodale MA, Mansfield RJW, eds. *Analysis of Visual Behavior*. Cambridge, MA: MIT Press; 1982:549-586.
83. Logie RH. *Visuo-Spatial Working Memory*. Hove, U.K.: Erlbaum; 1995.
84. Logie RH. Spatial and visual working memory: a mental workspace. *Psychology of Learning and Motivation: Advances and Theory: Cognitive Vision*. 2003;42:37-78.
85. Cowan N. Evolving conceptions of memory storage, selective attention, and their mutual constraints within the human information processing system. *Psychological Bulletin*. 1988;104:163-171.
86. Cowan N. *Attention and Memory: An Integrated Framework*. New York: Oxford University Press; 1995.
87. Cowan N. An embedded-processes model of working memory. In: Miyake A, Shah P, eds. *Models of Working Memory: Mechanisms of Active Maintenance and Executive Control*. Cambridge, U.K.: Cambridge University Press; 1999:62-101.
88. 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.
89. Oberauer K. Control of the contents of working memory—A comparison of two paradigms and two age groups. *Journal of Experimental Psychology: Learning, Memory, & Cognition*. 2005;31:714-728.
90. Oberauer K. The focus of attention in working memory—from metaphors to mechanisms. *Frontiers in Human Neuroscience*. 2013;7:doi:10.3389/fnhum.2013.00673.
91. van den Berg R, Awh E, Ma WJ. Factorial comparison of working memory models. *Psychological Review*. 2014;121:124-149.
92. Oberauer K, Kliegl R. A formal model of capacity limits in working memory. *Journal of Memory and Language*. 2006;55(4):601-626.

93. Gosmann J, Eliasmith C. A spiking neural model of the n-back task. In: Noelle DC, Dale R, Warlaumont AS, et al., eds. *Proceedings of the 37th Annual Meeting of the Cognitive Science Society*. Austin, TX: Cognitive Science Society; 2015:812-817.
94. Bays PM. Noise in neural populations accounts for errors in working memory. *Journal of Neuroscience*. 2014;34:3632-3645.
95. Edin F, Klingberg T, Johansson P, McNab F, Tegnér J, Compte A. Mechanism for top-down control of working memory capacity. *Proceedings of the National Academy of Sciences*. 2009;106:6802-6807.
96. Wei Z, Wang X-J, Wang D-H. From distributed resources to limited slots in multiple-item working memory: A spiking network model with normalization. *Journal of Neuroscience*. 2012(32).
97. Mongillo G, Barak O, Tsodyks M. Synaptic theory of working memory. *Science*. 2008;319:1543-1546.
98. Barak O, Tsodyks M. Working models of working memory. *Current Opinion in Neurobiology*. 2014;25:20-24.
99. Oberauer K, Lewandowsky S, Farrell S, Jarrold C, Greaves M. Modeling working memory: An interference model of complex span. *Psychonomic Bulletin & Review*. 2012;19:779-819.
100. Schneegans S, Bays PM. Neural architecture for feature binding in visual working memory. *Journal of Neuroscience*. 2017;37:3913-3925.
101. Awh E, Vogel EK. In: Kahana MJ, Wagner AD, eds. *Oxford Handbook of Memory*. in press.
102. Oberauer K, Farrell S, Jarrold C, Lewandowsky S. What limits working memory capacity? *Psychological Bulletin*. 2016;142:758-799.
103. Schweickert R, Boruff B. Short-term memory capacity: Magic number or magic spell? *Journal of Experimental Psychology: Learning, Memory, and Cognition*. 1986;12:419-425.
104. Cowan N, Blume CL, Saults JS. Attention to attributes and objects in working memory. *Journal of Experimental Psychology: Learning, Memory, and Cognition*. 2013;39:731-747.
105. Zhang W, Luck SJ. Discrete fixed-resolution representations in visual working memory. *Nature*. 2008;453:233-236.
106. Ma WJ, Husain M, Bays PM. Changing concepts of working memory. *Nature Neuroscience*. 2014;17:347-356.
107. Norman DA, Bobrow DG. On data-limited and resource-limited processes. *Cognitive Psychology*. 1975;7:44-64.
108. Lovett MC, Reder LM, Lebiere C. *Modeling individual differences in a digit working memory task*. Hillsdale: Erlbaum; 1997.
109. van den Berg R, Shin H, Chou W-C, George R, Ma WJ. Variability in encoding precision accounts for visual short-term memory limitations. *Proceedings of the National Academy of Sciences*. 2012;109:8780-8785.
110. Oberauer K, Lewandowsky S, Farrell S, Jarrold C, Greaves M. Modeling working memory: an interference model of complex span. *Psychonomic Bulletin & Review*. 2012;19:779-819.
111. Chase WG, Simon HA. Perception in Chess. *Cognitive Psychology*. 1973;4:55-81.
112. Gobet F, Lane PCR, Croker S, et al. Chunking mechanisms in human learning. *Trends in Cognitive Sciences*. 2001;5:236-243.
113. Hulme C, Roodenrys S, Schweickert R, Brown GDA, Martin S, Stuart G. Word-frequency effects on short-term memory tasks: Evidence for a redintegration process in immediate serial recall. *Journal of Experimental Psychology: Learning, Memory, and Cognition*. 1997;23:1217-1232.
114. Schweickert R. A multinomial processing tree model for degradation and redintegration in immediate recall. *Memory & Cognition*. 1993;21:168-173.
115. Romani C, McAlpine S, Martin RC. Concreteness effects in different tasks: Implications for models of short-term memory. *Quarterly Journal of Experimental Psychology*. 2008;61:292-323.

116. Jalbert A, Neath I, Surprenant AM. Does length or neighborhood size cause the word length effect. *Memory & Cognition*. 2011;39:1198-1210.
117. Thalmann M, Souza AS, Oberauer K. How does chunking help working memory? *Journal of Experimental Psychology: Learning, Memory and Cognition*. 2019;45:37-55.
118. Brown GDA, Neath I, Chater N. A temporal ratio model of memory. *Psychological Review*. 2007;114:539-576.
119. Sederberg PB, Howard MC, Kahana MJ. A context-based theory of recency and contiguity in free recall. *Psychological Review*. 2008;115:893-912.
120. Jeneson A, Squire L. Working memory, long-term memory, and medial temporal lobe function. *Learning & Memory*. 2012;19:15-25.
121. Watson PD, Voss JL, Warren DE, Tranel D, Cohen NJ. Spatial reconstruction by patients with hippocampal damage is dominated by relational memory errors. *Hippocampus*. 2013;23(7):570-580.
122. Pertzov Y, Miller TD, Gorgoraptis N, et al. Binding deficits in memory following medial temporal lobe damage in patients with voltage-gated potassium channel complex antibody-associated limbic encephalitis. *Brain*. 2013;136(8):2474-2485.
123. Papagno C, Shallice T. Introduction to impairments of short-term memory buffers: Do they exist? *Cortex*. 2019;112:1-4.
124. Wickens DD, Moody MJ, Dow R. The nature and timing of the retrieval process and of interference effects. *Journal of Experimental Psychology: General*. 1981;110:1-20.
125. Davelaar EJ, Goshen-Gottstein Y, Ashkenazi A, Haarmann HJ, Usher M. The demise of short-term memory revisited: empirical and computational investigation of recency effects. *Psychological Review*. 2005;112:3-42.
126. Beaudry O, Neath I, Surprenant AM, Tehan G. The focus of attention is similar to other memory systems rather than uniquely different. *frontiers in Human Neuroscience*. 2014;8. doi:10.3389/fnhum.2014.00056.
127. Norris D. Short-term memory and long-term memory are still different. *Psychological Bulletin*. 2017;143:992-1009.
128. Cowan N. Short-term memory based on activated long-term memory: A review in response to Norris (2017). *Psychological Bulletin*. 2019;145:822-847.
129. Norris D. Even an activated long-term memory system still needs a separate short-term store: A reply to Cowan (2019). *Psychological Bulletin*. 2019;145:848-853.
130. Farrell S. Temporal clustering and sequencing in short-term memory and episodic memory. *Psychological Review*. 2012;119:223-271.
131. Cowan N. *Working memory capacity*. New York: Psychology Press; 2005.
132. Oberauer K. Working memory and attention - a conceptual analysis and review. *Journal of Cognition*. 2019;2:1-23.
133. Case R, Kurland M, Goldberg J. Operational efficiency and the growth of short-term memory span. *Journal of Experimental Child Psychology*. 1982;33:386-404.
134. Just MA, Carpenter PA. A capacity theory of comprehension: Individual differences in working memory. *Psychological Review*. 1992;99:122-149.
135. Barrouillet P, Portrat S, Camos V. On the law relating processing to storage in working memory. *Psychological Review*. 2011;118:175-192.
136. Pashler H. Dual-task interference in simple tasks: Data and theory. *Psychological Bulletin*. 1994;116:220-244.
137. Jolicoeur P, Dell'Acqua R. The demonstration of short-term consolidation. *Cognitive Psychology*. 1998;36:138-202.
138. Vergauwe E, Camos V, Barrouillet P. The impact of storage on processing: How is information maintained in working memory? *Journal of Experimental Psychology: Learning, Memory, and Cognition*. 2014;40:1072-1095.

139. Thalmann M, Souza AS, Oberauer K. Revisiting the attentional demands of rehearsal in working-memory tasks. *Journal of Memory and Language*. 2019;105:1-18.
140. Oberauer K, Demmrich A, Mayr U, Kliegl R. Dissociating retention and access in working memory: An age-comparative study of mental arithmetic. *Memory & Cognition*. 2001;29(1):18-33.
141. Hazeltine E, Witfall T. Searching working memory for the source of dual-task costs. *Psychological Research*. 2011;75:466-475.
142. Klapp ST, Marshburn EA, Lester PT. Short-term memory does not involve the "working memory" of information processing: The demise of a common assumption. *Journal of Experimental Psychology: General*. 1983;112:240-264.
143. Tsubomi H, Fukuda K, Watanabe K, Vogel EK. Neural limits to representing objects still within view. *Journal of Neuroscience*. 2013;33:8257-8263.
144. Ester EF, Fukuda K, May LM, Vogel EK, Awh E. Evidence for a fixed capacity limit in attending multiple locations. *Cognitive, Affective, & Behavioral Neuroscience*. 2014;14:62-77.
145. Santangelo V, Macaluso E. The contribution of working memory to divided attention. *Human Brain Mapping*. 2013;34:158-175.
146. Souza AS, Oberauer K. The contributions of visual and central attention to visual working memory. *Attention, Perception & Psychophysics*. 2017;79:1897-1916.
147. Matsukura M, Vecera SP. Interference between object-based attention and object-based memory. *Psychonomic Bulletin & Review*. 2009;16:529-536.
148. Woodman GF, Chun MM. The role of working memory and long-term memory in visual search. *Visual Cognition*. 2006;14:808-830.
149. Wickens DD. Some characteristics of word encoding. *Memory & Cognition*. 1973;1:485-490.
150. Postle BR, Idzikowski C, Della Salla S, Logie RH, Baddeley AD. The selective disruption of spatial working memory by eye movements. *Q J Exp Psychol*. 2006;59:100-120.
151. Postle BR, Hamidi M. Nonvisual codes and nonvisual brain areas support visual working memory. *Cerebral Cortex*. 2007;17:2134-2142.
152. Postle BR, D'Esposito M, Corkin S. Effects of verbal and nonverbal interference on spatial and object visual working memory. *Memory & Cognition*. 2005.
153. Kane MJ, Conway ARA, Hambrick DZ, Engle RW. Variation in working memory capacity as variation in executive attention and control. In: Conway ARA, Jarrold C, Kane MJ, Miyake A, Towse JN, eds. *Variation in working memory*. New York: Oxford University Press; 2007:21-48.
154. Kane MJ, Bleckley MK, Conway ARA, Engle RW. A controlled-attention view of working-memory capacity. *Journal of Experimental Psychology: General*. 2001;130:169-183.
155. Baddeley AD. Exploring the central executive. *Quarterly Journal of Experimental Psychology*. 1996;49A:5-28.
156. Allen RJ, Baddeley AD, Hitch GJ. Is the binding of visual features in working memory resource-demanding? *Journal of Experimental Psychology: General*. 2006;135:298-313.
157. Lavie N. Distracted and confused?: Selective attention under load. *Trends in Cognitive Sciences*. 2005;9:75-82.
158. Rey-Mermet A, Gade M, Souza AS, von Bastian CC, Oberauer K. Is executive control related to working memory capacity and fluid intelligence? *Journal of Experimental Psychology: General*. 2019.
159. Rey-Mermet A, Gade M, Oberauer K. Should we stop thinking about inhibition? Searching for individual and age differences in inhibitory ability. *Journal of Experimental Psychology: Learning, Memory and Cognition*. 2018;44:501-526.
160. Kelley TA, Lavie N. Working Memory Load Modulates Distractor Competition in Primary Visual Cortex. *Cerebral Cortex*. 2011;21(3):659-665.
161. Lavie N, Hirst A, de Fockert JW, Viding E. Load theory of selective attention and cognitive control. *Journal of Experimental Psychology: General*. 2004;133:339-354.

162. SanMiguel I, Corral M-J, Escera C. When loading working memory reduces distraction: Behavioral and electrophysiological evidence from an auditory-visual distraction paradigm. *Journal of Cognitive Neuroscience*. 2008;20:1131-1145.
163. Scharinger C, Soutschek A, Schubert T, Gerjets P. When flanker meets the n-back: What EEG and pupil dilation data reveal about the interplay between the two central-executive working memory functions inhibition and updating. *Psychophysiology*. 2015;52(10):1293-1304.
164. Konstantinou N, Beal E, King J-R, Lavie N. Working memory load and distraction: dissociable effects of visual maintenance and cognitive control. *Attention, Perception & Psychophysics*. 2014;76:1985-1997.
165. Konstantinou N, Lavie N. Dissociable roles of different types of working memory load in visual detection. *Journal of Experimental Psychology: Human Perception and Performance*. 2013;39:919–924.
166. Park S, Kim M-S, Chun MM. Concurrent working memory load can facilitate selective attention: Evidence for specialized load. *JEP:HPP*. 2007;33(1062-1075).
167. Kim S-Y, Kim M-S, Chun MM. Concurrent working memory load can reduce distraction. *Proceedings of the National Academy of Sciences*. 2005;102:16524-16529.
168. Oberauer K, Hein L. Attention to information in working memory. *Current Directions in Psychological Science*. 2012;21:164-169.
169. Olivers CNL, Peters J, Houtkamp R, Roelfsema PR. Different states in visual working memory: when it guides attention and when it does not. *Trends in Cognitive Sciences*. 2011;15:327-334.
170. Soto D, Hodsoll J, Rotshtein P, Humphreys GW. Automatic guidance of attention from working memory. *Trends in Cognitive Sciences*. 2008;12:342-348.
171. Jacobsen CF. The functions of the frontal association areas in monkeys. *Comparative Psychology Monographs*. 1936;13:1-60.
172. Malmö RB. Interference factors in delayed response in monkey after removal of the frontal lobes. *Journal of Neurophysiology*. 1942;5:295-308.
173. Mishkin M. Effects of small frontal lesions on delayed alternation in monkeys. *Journal of Neurophysiology*. 1957;20:615-622.
174. Mishkin M, Pribram KH. Analysis of the effects of frontal lesions in the monkey. I. Variations of delayed alternation. *Journal of Comparative and Physiological Psychology*. 1955;48:492-495.
175. Pribram KH, Tubbs WE. Short-term memory, parsing, and the primate frontal cortex. *Science*. 1967;156:1765-1767.
176. Pribram KH, Ahumada A, Hartog J, Roos L. A progress report on the neurological processes disturbed by frontal lesions in primates. In: Warren JM, Akert K, eds. *The Frontal Granular Cortex and Behavior*. New York: McGraw-Hill Book Company; 1964:28-55.
177. Westbrook JA, Frank MJ. Dopamine and proximity in motivation and cognitive control. *Current Opinion in Behavioral Sciences*. 2018;22:28-34.
178. Rougier NP, Noelle DC, Braver TS, Cohen JD, O'Reilly RC. Prefrontal cortex and flexible cognitive control: Rules without symbols. *Proceedings of the National Academy of Sciences (USA)*. 2005;102:7338-7343.
179. Ghent L, Mishkin M, Teuber H-L. Short-term memory after frontal-lobe injury in man. *Journal of Comparative and Physiological Psychology*. 1962;5:705-709.
180. Milner B. Some effects of frontal lobectomy in man. In: Warren JM, Akert K, eds. *The Frontal Granular Cortex and Behavior*. New York: McGraw-Hill; 1964:313-334.
181. Chao LL, Knight RT. Human prefrontal lesions increase distractibility to irrelevant sensory inputs. *NeuroReport*. 1995;6:1605-1610.
182. Chao LL, Knight RT. Contribution of human prefrontal cortex to delay performance. *J Cog Neuroscience*. 1998;10:167-177.

183. Knight RT, Hillyard SA, Woods DL, Neville HJ. The effects of frontal cortex lesions on event-related potentials during auditory selective attention. *Electroencephalography and Clinical Neurophysiology*. 1981;52:571-582.
184. Hazy TE, Frank MJ, O'Reilly RC. Towards an executive without a homunculus: Computational models of the prefrontal cortex/basal ganglia system. *Philosophical Transactions of the Royal Society of London Series B, Biological Sciences*. 2007;362:1601–1613.
185. Badre D, Nee DE. Frontal cortex and the hierarchical control of behavior. *Trends in Cognitive Sciences*. 2018;22:170-188.
186. Christophel TB, Klink PC, Spitzer B, Roelfsema PR, Haynes J-D. The distributed nature of working memory. *Trends in Cognitive Sciences*. 2017;21:111-124.
187. Leavitt ML, Mendoza-Halliday D, Martinez-Trujillo JC. Sustained activity encoding working memories: not fully distributed. *Trends in Neurosciences*. 2018;40:328-346.
188. Constantinides C, Funahashi S, Lee D, et al. Persistent spiking activity underlies working memory. *J Neurosci*. 2018;38:7020-7028.
189. Jacobsen CF, Nissen HW. Studies of cerebral function in primates. IV. The effects of frontal lobe lesions on the delayed alternation habit in monkeys. *Journal of Comparative Psychology*. 1937;23:101-112.
190. Fuster JM, Alexander GE. Neuron activity related to short-term memory. *Science*. 1971;173:652-654.
191. Kubota K, Niki H. Prefrontal cortical unit activity and delayed alternation performance in monkeys. *J Neurophysiol*. 1971;34(3):337-347.
192. Postle BR. Neural bases of the short-term retention of visual information. In: Jolicoeur P, LeFebvre C, Martinez-Trujillo J, eds. *Mechanisms of Sensory Working Memory: Attention & Performance XXV*. London, U.K.: Academic Press; 2015:43-58.
193. Goldman-Rakic PS. Circuitry of the prefrontal cortex and the regulation of behavior by representational memory. In: Mountcastle VB, Plum F, Geiger SR, eds. *Handbook of Neurobiology*. Bethesda: American Physiological Society; 1987:373-417.
194. Funahashi S, Bruce CJ, Goldman-Rakic PS. Mnemonic coding of visual space in the monkey's dorsolateral prefrontal cortex. *Journal of Neurophysiology*. 1989;61:331-349.
195. Funahashi S, Bruce CJ, Goldman-Rakic PS. Visuospatial coding in primate prefrontal neurons revealed by oculomotor paradigms. *Journal of Neurophysiology*. 1990;63:814-831.
196. Funahashi S, Chafee MV, Goldman-Rakic PS. Prefrontal neuronal activity in rhesus monkeys performing a delayed anti-saccade task. *Nature*. 1993;365:753-756.
197. Wilson FAW, O'Scalaidhe SP, Goldman-Rakic PS. Dissociation of object and spatial processing domains in primate prefrontal cortex. *Science*. 1993;260:1955-1958.
198. Funahashi S, Bruce C, Goldman-Rakic P. Dorsolateral prefrontal lesions and oculomotor delayed-response performance: Evidence for mnemonic "scotomas". *J Neurosci*. 1993;13:1479-1497.
199. Davachi L, Romanski LM, Chafee MV, Goldman-Rakic PS. Domain specificity in cognitive systems. In: Gazzaniga MS, ed. *The Cognitive Neurosciences III*. Cambridge, MA: The MIT Press; 2004:665-678.
200. Goldman-Rakic PS, Leung H-C. Functional architecture of the dorsolateral prefrontal cortex in monkeys and humans. In: Stuss DT, Knight RT, eds. *Principles of Frontal Lobe Function*. Oxford, U.K.: Oxford University Press; 2002:85-95.
201. Rainer G, Asaad WF, Miller EK. Memory fields of neurons in the primate prefrontal cortex. *Proceedings of the National Academy of Sciences (USA)*. 1998;95:15008-15013.
202. Rao SC, Rainer G, Miller EK. Integration of what and where in the primate prefrontal cortex. *Science*. 1997;276:821-824.
203. Lebedev MA, Messinger A, Kralik JD, Wise SP. Representation of attended versus remembered locations in prefrontal cortex. *PLoS Biology*. 2004;2:1919-1935.

204. Tsujimoto S, Postle BR. The prefrontal cortex and delay tasks: a reconsideration of the "mnemonic scotoma". *J Cog Neuroscience*. 2012;24:627-635.
205. Funahashi S. Functions of delay-period activity in the prefrontal cortex and mnemonic scotomas revisited. *Frontiers in systems neuroscience*. 2015;9.
206. Riley MR, Constantinidis C. The role of prefrontal persistent activity in working memory. *Frontiers in systems neuroscience*. 2016.
207. 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.
208. Jonikaitis D, Moore T. The interdependence of attention, working memory and gaze control: behavior and neural circuitry. *Current Opinion in Psychology*. 2019;29:126-134.
209. LaRocque JJ, Lewis-Peacock JA, Postle BR. Multiple neural states of representation in short-term memory? It's a matter of attention. *Frontiers in Human Neuroscience*. 2014;8:doi:10.3389/fnhum.2014.00005.
210. Moore T, Fallah M. Microstimulation of the frontal eye field and its effects on covert attention. *Journal of Neurophysiology*. 2004;91:152-162.
211. Moore T, Fallah M. Control of eye movements and spatial attention. *Proceedings of the National Academy of Science (USA)*. 2001;98:1273-1276.
212. Cavanaugh J, Wurtz RH. Subcortical modulation of attention counters change blindness. *J Neurosci*. 2004;24:11236-11243.
213. Muller JR, Philiastides MG, Newsome WT. Microstimulation of the superior colliculus focuses attention without moving the eyes. *Proceedings of the National Academy of Science (USA)*. 2005;102:524-529.
214. Moore T, Armstrong KM. Selective gating of visual signals by microstimulation of frontal cortex. *Nature*. 2003;421:370-373.
215. Armstrong KM, Chang MH, Moore T. Selection and maintenance of spatial information by frontal eye field neurons. *J Neurosci*. 2009;29:15621-15629.
216. Mirpour K, Bolandnazar S, Bisley JW. Neurons in FEF keep track of items that have been previously fixated in free viewing visual search. *J Neurosci*. 2019;39:2114-2124.
217. Clark KL, Noudoost B, Moore T. Persistent spatial information in the FEF during object-based short-term memory does not contribute to task performance. *J Cog Neuroscience*. 2014;26:1292-1299.
218. Jerde T, Merriam EP, Riggall AC, Hedges JH, Curtis CE. Prioritized Maps of Space in Human Frontoparietal Cortex. *The Journal of Neuroscience*. 2012;32:17382-17390.
219. Mackey W, Devinsky O, Doyle W, Meager M, Curtis CE. Human dorsolateral prefrontal cortex is not necessary for spatial working memory. *J Neurosci*. 2016;36:2847-2856.
220. Hamidi M, Tononi G, Postle BR. Evaluating frontal and parietal contributions to spatial working memory with repetitive transcranial magnetic stimulation. *Brain Research*. 2008;1230:202-210.
221. Baldauf D, Desimone R. Neural mechanisms of object-based attention. *Science*. 2014;344:424-427.
222. Bichot NP, Heard MT, DeGennaro EM, Desimone R. A source for feature-based attention in prefrontal cortex. *Neuron*. 2015;88:832-844.
223. Compte A, Brunel N, Goldman-Rakic PS, Wang X-J. Synaptic mechanisms and network dynamics underlying spatial working memory in a cortical network model. *Cerebral Cortex*. 2000;10:910-923.
224. Machens CK, Romo R, Brody CD. Flexible control of mutual inhibition: a neural model of two-interval discrimination. *Science*. 2005;307:1121-1124.
225. Murray JD, Bernacchia A, Roy NA, Constantinidis C, Romo RX, Wang X-J. Stable population coding for working memory coexists with heterogeneous neural dynamics in prefrontal cortex. *Proceedings of the National Academy of Sciences (USA)*. 2017;114:394-399.

226. Wang X-J. Synaptic reverberation underlying mnemonic persistent activity. *Trends in Neurosciences*. 2001;24:455-463.
227. Mendoza-Halliday D, Torres S, Martinez-Trujillo JC. Sharp emergence of feature-selective sustained activity along the dorsal visual pathway. *Nature Neuroscience*. 2014;17:1255-1262.
228. Mendoza-Halliday D, Martinez-Trujillo JC. Neuronal population coding of perceived and memorized visual features in the lateral prefrontal cortex. *Nature Communications*. 2017;8.
229. Hebb DO. *Organization of Behavior*. New York, NY: John Wiley & Sons, Inc.; 1949.
230. Vogel EK, Machizawa MG. Neural activity predicts individual differences in visual working memory capacity. *Nature*. 2004;428:748-751.
231. Todd JJ, Marois R. Capacity limit of visual short-term memory in human posterior parietal cortex. *Nature*. 2004;428:751-754.
232. Todd JJ, Marois R. Posterior parietal cortex activity predicts individual differences in visual short-term memory capacity. *Cognitive, Affective, & Behavioral Neuroscience*. 2005;5:144-155.
233. Xu Y, Chun MM. Dissociable neural mechanisms supporting visual short-term memory for objects. *Nature*. 2006;440:91-95.
234. Norman KA, Polyn SM, Detre GJ, Haxby JV. Beyond mind-reading: multi-voxel pattern analysis of fMRI data. *Trends in Cognitive Sciences*. 2006;10:424-430.
235. Haxby JV, Gobbini MI, Furey ML, Ishai A, Schouten JL, Pietrini P. Distributed and overlapping representations of faces and objects in ventral temporal cortex. *Science*. 2001;293:2425-2430.
236. Pereira F, Mitchell T, Botvinick MM. Machine learning classifiers and fMRI: a tutorial overview. *NeuroImage*. 2009;45:S199-S209.
237. Riggall AC, Postle BR. The relationship between working memory storage and elevated activity as measured with functional magnetic resonance imaging. *The Journal of Neuroscience*. 2012;32:12990-12998.
238. Lewis-Peacock JA, Postle BR. Temporary activation of long-term memory supports working memory. *The Journal of Neuroscience*. 2008;28:8765-8771.
239. Harrison SA, Tong F. Decoding reveals the contents of visual working memory in early visual areas. *Nature*. 2009;458:632-635.
240. Serences JT, Ester EF, Vogel EK, Awh E. Stimulus-specific delay activity in human primary visual cortex. *Psychological Sci*. 2009;20:207-214.
241. Sprague TC, Serences JT. Attention modulates spatial priority maps in the human occipital, parietal and frontal cortices. *Nature Neuroscience*. 2013;16:1879-1887.
242. Sprague TC, Ester EF, Serences JT. Reconstructions of information in visual spatial working memory degrade with memory load. *Current Biology*. 2014;24:2174-2180.
243. Emrich SM, Riggall AC, Larocque JJ, Postle BR. Distributed patterns of activity in sensory cortex reflect the precision of multiple items maintained in visual short-term memory. *J Neurosci*. 2013;33:6516-6523.
244. Gosseries O, Yu Q, LaRocque JJ, et al. Parieto-occipital interactions underlying control- and representation-related processes in working memory for nonspatial visual features. *J Neurosci*. 2018;38:4357-4366.
245. Ester EF, Anderson DE, Serences JT, Awh E. A neural measure of precision in visual working memory. *Journal of Cognitive Neuroscience*. 2013;25(5):754-761.
246. van Ede F, Chekroud SR, Stokes MG, Nobre AC. Concurrent visual and motor selection during visual working memory guided action. *Nature Neuroscience*. 2019;22:477-483.
247. Serences JT. Neural mechanisms of information storage in visual short-term memory. *Vision Research*. 2016;128:53-67.
248. D'Esposito M, Postle BR. The cognitive neuroscience of working memory. *Annual Review of Psychology*. 2015;66:115-142.

249. Kar K, Kubilius J, Schmidt K, Issa EB, DiCarlo JJ. Evidence that recurrent circuits are critical to the ventral stream's execution of core object recognition behavior. *Nature Neuroscience*. 2019;22:974-983.
250. Cudeiro J, Sillito AM. Looking back: corticothalamic feedback and early visual processing. *Trends in Neuroscience*. 2006;29(6):298-306.
251. Sillito AM, Cudeiro J, Jones HE. Always returning: feedback and sensory processing in visual cortex and thalamus. *Trends in Neuroscience*. 2006;29(6):307-316.
252. Lamme VAF. The neurophysiology of figure-ground segregation in primary visual cortex. *J Neurosci*. 1995;15:1605-1615.
253. Cai Y, Sheldon AD, Postle BR. The influence of storage capacity versus control in visual working memory capacity. *Annual Meeting of the Cognitive Neuroscience Society*. 2018;March 2018.
254. van Kerkoerle T, Self MW, Roelfsema PR. Layer-specificity in the effects of attention and working memory on activity in primary visual cortex. *Nature Communications*. 2017;8:13804.
255. Rademaker RL, Chunharas C, Serences JT. Coexisting representations of sensory and mnemonic information in human visual cortex. *Nature Neuroscience*. 2019;22:1336-1344.
256. Christophel TB, Hebart MN, Haynes J-D. Decoding the contents of visual short-term memory from human visual and parietal cortex. *The Journal of Neuroscience*. 2012;32:2983-12989.
257. Christophel TB, Iamshchinina P, Yan C, Allefeld C, Haynes J-D. Cortical specialization for attended versus unattended working memory. *Nature Neuroscience*. 2018;21:494-496.
258. Christophel TB, Haynes JD. Decoding complex flow-field patterns in visual working memory. *NeuroImage*. 2014;91:43-51.
259. Bettencourt KC, Xu Y. Decoding the content of visual short-term memory under distraction in occipital and parietal areas. *Nature Neuroscience*. 2016;19:150-157.
260. Ester EF, Sprague TC, Serences JT. Parietal and frontal cortex encode stimulus-specific mnemonic representations during visual working memory. *Neuron*. 2015;87:893-905.
261. Cai Y, Sheldon AD, Yu Q, Postle BR. Overlapping and distinct contributions of stimulus location and of spatial context to nonspatial visual short-term memory. *Journal of Neurophysiology*. 2019;121:1222-1231.
262. Yu Q, Shim WM. Occipital, parietal, and frontal cortices selectively maintain task-relevant features of multi-feature objects in visual working memory. *NeuroImage*. 2017;157:97-107.
263. Xu Y. Reevaluating the sensory account of visual working memory storage. *Trends in Cognitive Sciences*. 2017;27:794-815.
264. Xu Y. Sensory cortex is nonessential in working memory storage. *Trends in Cognitive Sciences*. 2018;22:192-193.
265. Scimeca JM, Kiyonaga A, M. DE. Reaffirming the sensory recruitment account of working memory. *Trends in Cognitive Sciences*. 2018;22:190-192.
266. Gayet S, Paffen CLE, Van der Stigchel S. Visual working memory storage recruits sensory processing areas. *Trends in Cognitive Sciences*. 2018;22:189-190.
267. Carlisle NB, Arita JT, Pardo D, Woodman GF. Attentional templates in visual working memory. *The Journal of Neuroscience*. 2011;31:9315-9322.
268. Reinhart RM, Woodman GF. High stakes trigger the use of multiple memories to enhance the control of attention. *Cerebral Cortex*. 2014;24:2022-2035.
269. Panichello MF, DePasquale B, Pillow JW, Buschman TJ. Error-correcting dynamics in visual working memory. *Nature Communications*. 2019;10:3366.
270. Cai Y, Yu Q, Sheldon AD, Postle BR. The role of location-context binding in nonspatial visual working memory. *bioRxiv*. unpublished.
271. Yu Q, Panichello MF, Postle BR, Buschman TJ. Persistent neural activity in parietal cortex tracks attractor dynamics in visual working memory. *Poster presented at the Annual Meeting of the Cognitive Neuroscience Society*. 2019:March 2019, San Francisco, CA.

272. Fiebig F, Lansner A. A spiking working memory model based on Hebbian short-term potentiation. *The Journal of Neuroscience*. 2017;37:83-96.
273. Lundqvist M, Herman P, Lansner A. Theta and gamma power increases and alpha/beta power decreases with memory load in an attractor network model. *Journal of Cognitive Neuroscience*. 2011;23:3008-3020.
274. Lundqvist M, Compte A, Lansner A. Bistable, irregular firing and population oscillations in a modular attractor memory network. *PLoS Computational Biology*. 2010;6:e1000803.
275. Lundqvist M, Rose J, Herman P, Brincat SL, Buschman TJ, Miller EK. Gamma and beta bursts underlie working memory. *Neuron*. 2016;90:1-13.
276. Lundqvist M, Herman P, Miller EK. Working Memory: Delay Activity, Yes! Persistent Activity? Maybe Not. *J Neurosci*. 2018;38:7013-7019.
277. Constantinidis C, Funahashi S, Lee D, et al. Persistent Spiking Activity Underlies Working Memory. *J Neurosci*. 2018;38:7020-7028.
278. Miller EK, Lundqvist M, Bastos AM. Working Memory 2.0. *Neuron*. 2018;100:463-475.
279. Itskov V, Hansel D, Tsodyks M. Short-Term Facilitation may Stabilize Parametric Working Memory Trace. *Frontiers in computational neuroscience*. 2011;5:40.
280. Hayden BY, Gallant JL. Working memory and decision processes in visual area V4. *Frontiers in Neuroscience*. 2013;7:doi: 10.3389/fnins.2013.00018.
281. Sugase-Miyamoto Y, Liu Z, Wiener MC, Optican LM, Richmond BJ. Short-term memory trace in rapidly adapting synapses of inferior temporal cortex. *PLoS Computational Biology*. 2008;4(5):e1000073.
282. Stokes MG. 'Activity-silent' working memory in prefrontal cortex: a dynamic coding framework. *Trends in Cognitive Sciences*. 2015;19:394-405.
283. Wolff MJ, Ding J, Myers NE, Stokes MG. Revealing hidden states in visual working memory using electroencephalography. *Frontiers in systems neuroscience*. 2015;9.
284. Wolff MJ, Jochim J, Akyürek EG, Stokes MG. Dynamic hidden states underlying working-memory-guided behavior. *Nature Neuroscience*. 2017.
285. Lewis-Peacock JA, Drysdale AT, Oberauer K, Postle BR. Neural evidence for a distinction between short-term memory and the focus of attention. *Journal of Cognitive Neuroscience*. 2012;24:61-79.
286. LaRocque JJ, Lewis-Peacock JA, Drysdale A, Oberauer K, Postle BR. Decoding attended information in short-term memory: An EEG study. *J Cog Neuroscience*. 2013;25:127-142.
287. Rose N, Larocque JJ, Riggall AC, et al. Reactivation of latent working memories with transcranial magnetic stimulation. *Science*. 2016;354:1136-1139.
288. LaRocque JJ, Riggall AC, Emrich SM, Postle BR. Within-category decoding of information in different states in short-term memory. *Cerebral Cortex*. 2017;17:4881-4890.
289. Masse NY, Yang GR, Song HF, Wang X-J, Freedman DJ. Circuit mechanisms for the maintenance and manipulation of information in working memory. *Nature Neuroscience*. 2019;22:1159-1167.
290. Manohar SG, Zokaei N, Fallon SJ, Vogels TP, Husain M. Neural mechanisms of attending to items in working memory. *Neuroscience and Biobehavioral Reviews*. 2019;101:1-12.
291. Orhan AE, Ma WJ. A diverse range of factors affect the nature of neural representations underlying short-term memory. *Nature Neuroscience*. 2019;22:275-283.
292. Spaak E, Watanabe K, Funahashi S, Stokes MG. Stable and dynamic coding for working memory in primate prefrontal cortex. *The Journal of Neuroscience*. 2017;37:6503-6516.
293. Masse NY, Hodnefield JM, Freedman DJ. Mnemonic encoding and cortical organization in parietal and prefrontal cortices. *J Neurosci*. 2017;37:6098-6112.
294. Sarma A, Masse NY, Wang X-J, Freedman DJ. Task-specific versus generalized mnemonic representations in parietal and prefrontal cortices. *Nature Neuroscience*. 2016;19:143-149.
295. Zaksas D, Pasternak T. Directional signals in the prefrontal cortex and in area MT during a working memory for visual motion task. *The Journal of Neuroscience*. 2006;26:11726-11742.

296. Bisley JW, Mirpour K. The neural instantiation of a priority map. *Current Opinion in Psychology*. 2019;29:108-112.
297. Monsell S. Recency, immediate recognition memory, and reaction time. *Cognitive Psychology*. 1978;10:465-501.
298. Jonides J, Smith EE, Marshuetz C, Koeppel RA, Reuter-Lorenz PA. Inhibition of verbal working memory revealed by brain activation. *Proceedings of the National Academy of Sciences*. 1998;95:8410-8413.
299. Thompson-Schill S, Jonides J, Marshuetz C, et al. Effects of frontal lobe damage on interference effects in working memory. *Cognitive, Affective, and Behavioral Neuroscience*. 2002;2:109-120.
300. D'Esposito M, Postle BR, Jonides J, Smith EE. The neural substrate and temporal dynamics of interference effects in working memory as revealed by event-related functional MRI. *Proceedings of the National Academy of Sciences, USA*. 1999;96:7514-7519.
301. Feredoes E, Tononi G, Postle BR. Direct evidence for a prefrontal contribution to the control of proactive interference in verbal working memory. *Proceedings of the National Academy of Science (USA)*. 2006;103:19530-19534.
302. Feredoes E, Tononi G, Postle BR. Prefrontal control of familiarity vs. recollection in working memory. *Annual Meeting of the Cognitive Neuroscience Society*. 2007:Program # F 90.
303. Warrington EK, Shallice T. The selective impairment of auditory verbal short-term memory. *Brain*. 1969;92:885-896.
304. Logie RH. Covering sources of evidence and theory integration in working memory: A commentary on Morey, Rhodes, and Cowan (2019). *Cortex*. 2019;112:162-171.
305. Jonin P-Y, Calia C, Muratot S, et al. Refining understanding of working memory buffers through the construct of binding: Evidence from a single case informs theory and clinical practice. *Cortex*. 2019;112:37-57.
306. Tree JJ, Playfoot D. How to get by with half a loop — An investigation of visual and auditory codes in a case of impaired phonological short-term memory (pSTM). *Cortex*. 2019;112:23-36.
307. Martin RC, Schnur TT. Independent contributions of semantic and phonological working memory to spontaneous speech in acute stroke. *Cortex*. 2019;112:58-68.
308. Shallice T, Papagno C. Impairments of auditory-verbal short-term memory: Do selective deficits of the input phonological buffer exist? *Cortex*. 2019;112:107-121.
309. Hanley JR, Young AW. ELD revisited: A second look at a neuropsychological impairment of working memory affecting retention of visuo-spatial material. *Cortex*. 2019;112:172-179.
310. Morey CC. Working memory theory remains stuck: Reply to Hanley and Young. *Cortex*. 2019;112:180-181.
311. Saalman YB, Pinsk MA, Wang L, Li X, Kastner S. The pulvinar regulates information transmission between cortical areas based on attention demands. *Science*. 2012;337:753-756.
312. Moore T, Zirnsak M. Neural mechanisms of selective visual attention. *Annual Review of Psychology*. 2017;68:47-72.
313. Desimone R, Duncan J. Neural mechanisms of selective visual attention. *Annual Review of Neuroscience*. 1995;18:193-222.
314. Sheldon AD, Saad E, Sahan MI, et al. Attention biases competition for visual representation via enhancement of targets and inhibition of nontargets. unpublished.
315. Chatham CH, Badre D. Multiple gates on working memory. *Current Opinion in Behavioral Sciences*. 2015;1:23-31.
316. Desrochers TM, Burk DC, Badre D, Sheinberg DL. The monitoring and control of task sequences in human and non-human primates. *Frontiers in Systems Neuroscience*. 2015;9.
317. Lewis-Peacock JA, Kessler Y, Oberauer K. The removal of information from working memory. *Annals of the New York Academy of Science*. 2018;1424:33-44.

318. Myers NE, Stokes MG, Nobre AC. Prioritizing information during working memory: Beyond sustained internal attention. *Trends in Cognitive Sciences*. 2017;21:449-461.

Figure 1: Schematics of three theoretical frameworks of working memory: A: Multi-component model, with a central executive (CE) controlling three stores, the episodic buffer, the visual-spatial sketch pad (VSSP), and the phonological loop (Ph. Loop), the contents of which are re-activated through articulatory rehearsal. B: Embedded-process model: In a network of long-term memory representation a subset (grey nodes) forms the activated part of long-term memory. A limited number of 3-4 representations is in the focus of attention (thick black oval). C: Embedded-component model: A subset of the activated part of the long-term memory network is temporarily bound (thick broken lines) to a context (here: a location in a mental space, depicted as a spatial frame) and thereby related to each other (thick grey lines); this subset constitutes the direct-access region. Within it, one element and its context are selected into the focus of attention (thick black oval).

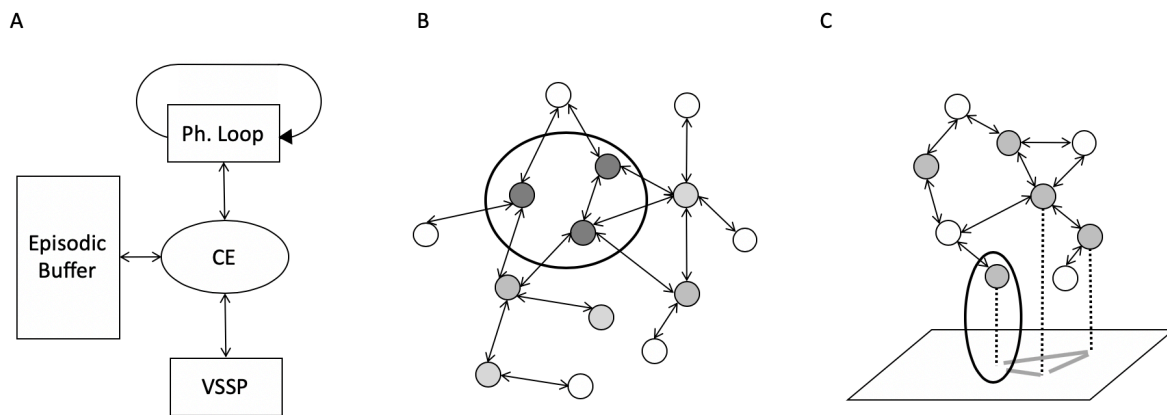


Figure 2: A: Architecture of the bump-attractor model⁹⁶ applied to working memory for orientations: Model neurons are arranged in a ring, ordered by the preferred orientation of their tuning functions; nearby neurons have excitatory connections (arrow heads, depicted outside the ring); distant neurons have inhibitory connections (nobs, inside the ring). Shading of the units reflects their activation (grey = baseline, black = above baseline; white = below baseline). B: Bump-attractor model after encoding a left-pointing arrow: An activation bump is created with a peak at the neuron preferentially coding left-ward orientations. C: Architecture of 2-layer neural network for binding contents to contexts. D: State of the 2-layer network after encoding a stimulus: The stimulus is represented as a pattern of activation in the content layer; its context as an activation pattern in the context layer. Rapid Hebbian learning updates the connection weights, strengthening those that are active with the same polarity.

