

# 1 Adaptive chunking improves 2 effective working memory capacity 3 in a prefrontal cortex and basal 4 ganglia circuit

5 Aneri V Soni<sup>1, §</sup> and Michael J. Frank<sup>1, ¶</sup>

\*For correspondence:

[aneri\\_soni@brown.edu](mailto:aneri_soni@brown.edu);  
[michael\\_frank@brown.edu](mailto:michael_frank@brown.edu)

6 <sup>1</sup>Brown University

7 **Present address:** <sup>§</sup>Neuroscience,  
Brown University, RI, USA; <sup>¶</sup>Dept of  
Cognitive, Linguistic, and  
Psychological Science, Brown  
University, RI, USA

8 **Abstract** How and why is working memory (WM) capacity limited? Traditional cognitive  
9 accounts focus either on limitations on the number or items that can be stored (slots models), or  
10 loss of precision with increasing load (resource models). Here we show that a neural network  
11 model of prefrontal cortex and basal ganglia can learn to reuse the same prefrontal populations  
12 to store multiple items, leading to resource-like constraints within a slot-like system, and inducing  
13 a trade-off between quantity and precision of information. Such “chunking” strategies are  
14 adapted as a function of reinforcement learning and WM task demands, mimicking human  
15 performance and normative models. Moreover, adaptive performance requires a dynamic range  
16 of dopaminergic signals to adjust striatal gating policies, providing a new interpretation of WM  
17 difficulties in patient populations such as Parkinson’s disease, ADHD and schizophrenia. These  
18 simulations also suggest a computational rather than anatomical limit to WM capacity.

## 20 Introduction

21 It has long been appreciated that working memory (WM) capacity is limited, but despite many  
22 decades of research, the nature of these limitations remains controversial. For example, early  
23 work by Miller (1951) famously posited that WM is limited to roughly 7 items, but he also stated that  
24 the precise limitation might vary depending on information quantity of the relevant memoranda.  
25 More recently, a vigorous debate over the last two decades has divided roughly into two schools of  
26 thought. The “slots” theory argues that WM is limited to a specific number of items. According to  
27 this account, each item is stored in a discrete slot and encoded with high precision, leading to low  
28 or zero errors when that item is probed at recall. When the number of items to be remembered  
29 exceeds the capacity (roughly 4 items; *Cowan (2008)*), some items will be forgotten, and therefore  
30 participants resort to guessing (*Zhang and Luck, 2008; Fukuda et al., 2010; Luck and Vogel, 2013*).  
31 In contrast, the “resource” theory argues that people can store an arbitrarily large number of items  
32 in WM, with no inherent item limit, but that each item competes for a shared pool of resources.  
33 As a result, the precision of each memoranda goes down with each added item to be recalled.  
34 In the limit, when many items are presented, each one is recalled with low precision, which can  
35 masquerade as guessing (*Bays et al., 2009; Ma et al., 2014*).

36 Critically, regardless of the distinction between discrete and continuous resources, the mea-  
37 sured WM “capacity” from experimental data is not fixed. For example, individual differences in  
38 WM capacity are largely determined not by the raw number of items one can store but rather one’s  
39 ability to filter out distracting items (*Vogel et al., 2005; McNab and Klingberg, 2008; Astle et al., 2014*;

40 *Feldmann-Wüstefeld and Vogel, 2019*). More generally, one can leverage various (potentially un-  
41 conscious) memory strategies to improve “effective capacity”, leading to experimentally observed  
42 capacity measurements that fluctuate depending on stimulus complexity, sensory modality, and  
43 experience with the stimuli (*Pusch et al., 2023*). Thus a key but often overlooked aspect of WM  
44 lies in the efficient *management* of access to and from working memory. Taking into account this  
45 management may also provide a mechanism for understanding WM not only in terms of mainte-  
46 nance, but also how gating strategies may be used for manipulation of information – i.e., “working  
47 with memory” (*Moscovitch and Winocur, 2009*). In other words, **effective capacity** encompasses  
48 gating items into/out of WM as well as the maintenance of items.

49 For example, recent theoretical and empirical work suggested that information stored in WM  
50 can be partitioned into discrete representations, but that similar items could be stored in a shared  
51 partition, in effect chunking them together (*Nassar et al., 2018*). Intuitively, it is simpler to remem-  
52 ber that one needs to purchase dairy items, bread, and fruit rather than to remember to buy milk,  
53 cheese, bread, oranges, and bananas. This active chunking strategy serves as a lossy information  
54 compression mechanism: it frees up space for other items to be stored and recalled. This happens  
55 at the cost of reducing precision for the chunked items. Experimental evidence provided support  
56 for such a model over alternatives, and provided a mechanism to explain why precision can be  
57 variable across trials as posited by previous resources models (*Berg et al., 2012*). Moreover, (*Nas-  
58 sar et al., 2018*) showed that the optimal chunking criterion (i.e., how similar two items need to  
59 be to merit storing as a single chunk) varies systematically with set size (number of items to be  
60 remembered) and can be acquired via reinforcement learning (RL). In line with this account, they  
61 reported evidence that humans adapted chunking on a trial by trial basis as a function of reward  
62 feedback in their experiment. They also performed a meta-analysis of other visual working mem-  
63 ory (VWM) datasets showed that optimal performance was associated with more chunking with  
64 increasing set size. Thus, at the cost of small errors (due to loss in precision), normative models  
65 and humans have overall better recall and performance when employing this chunking method.  
66 This and related theories (*Brady et al., 2011; Brady and Alvarez, 2015; van den Berg et al., 2014;  
67 Wei et al., 2012; Swan and Wyble, 2014; Nassar et al., 2018*) may reconcile differences between the  
68 slots and resources theories.

69 While these normative and algorithmic models are consistent with experimental data, it is un-  
70 clear how *biological* neural networks could perform such adaptive and flexible chunking. If a bio-  
71 logical neural network exhibits the same mechanism, we can make more clear predictions about  
72 human WM as well. The dominant neural model of visual working memory is the ring attractor  
73 model of prefrontal cortex (PFC), whereby multiple items can be maintained via persistent activity  
74 in attractor states (*Edin et al., 2009; Wei et al., 2012; Nassar et al., 2018*). In these models, nearby  
75 attractors coding for overlapping visual stimuli can collide, leading to a form of chunking and loss  
76 of precision, and where some items are forgotten due to lateral inhibition (*Wei et al., 2012; Almeida  
77 et al., 2015; Nassar et al., 2018*). While these models have successfully accounted for a range of  
78 data, by modeling only the PFC (or a single cortical population), they have limitations. Firstly, these  
79 models cannot determine whether or not an item should be stored. In other words, unlike hu-  
80 mans (*Vogel et al., 2005; McNab and Klingberg, 2008*), they cannot improve effective capacity by  
81 filtering content to only include relevant information. Secondly, any chunking that occurs in these  
82 models is obligatory – determined only by how overlapping the neural populations are and hence  
83 whether attractors will collide. Thus, chunking can’t be adapted with task demands as required  
84 by normative models and human data (*Nassar et al., 2018*). Finally, during recall, these network  
85 models cannot select a specific item from memory based on a probe (accuracy in these models is  
86 considered high as long as the relevant stimulus is encoded somewhere in the pool of neurons;  
87 (*Edin et al., 2009; Wei et al., 2012; Almeida et al., 2015*)). In other words, these models have no  
88 way of manipulating or accessing the contents of the WM.

89 This selective access and management of information in WM requires the brain to 1) solve the  
90 variable binding problem and 2) create a role- addressable memory. **Variable binding** refers to

91 the ability to link attributes of an object together (for instance a person has a name, face, location  
92 etc.) In WM, humans have to temporarily bind a given memorandum with a given slot in memory,  
93 often in a role-addressable manner. For example, they might need to recall the color of an object  
94 based on its shape. Indeed, limitations in WM capacity have been linked to difficulty in binding  
95 items to their respective slots in memory (**Oberauer, 2013; Oberauer and Lin, 2017**). Moreover, the  
96 selective updating of information in these slots is thought to produce a recency bias ubiquitously  
97 observed in human WM (**Oberauer et al., 2012**).

98 Another complementary line of biologically-inspired neural network models addresses how in-  
99 teractions between basal ganglia (BG), thalamus, and PFC support independent updating of sep-  
100 arable PFC “stripes” (anatomical clusters of interconnected neurons that are isolated from other  
101 stripes; (**Levitt et al., 1993; Pucak et al., 1996; Frank et al., 2001**). These prefrontal cortex basal  
102 ganglia working memory (PBWM) models focus on the aforementioned “management” and vari-  
103 able binding problem. They simulate the brain’s decision whether to encode a sensory item in  
104 WM (“selective input gating”). They also simulate which item (of those stored in WM) should be  
105 accessed (“output gating”) for reporting or subsequent processing (**O’Reilly and Frank, 2006; Hazy  
106 et al., 2007; Krueger and Dayan, 2009; Stocco et al., 2010; Frank and Badre, 2012; Kriete et al.,  
107 2013**). The combination of input and output gating decisions that are made can be summarized  
108 as the *gating policy*. Via dopaminergic reinforcement learning signaling in the BG, the networks  
109 learn an effective gating policy for a given WM task (**Frank and Badre, 2012**). This policy includes  
110 (i) whether or not to store an item (i.e., if it is task-relevant or distracting), (ii) if relevant, in which  
111 population of PFC neurons to store it, and (iii) which population of PFC neurons should be gated  
112 out during recall or action selection. As such, PBWM networks can perform complex tasks that  
113 require keeping track of sequences of events across multiple trials while also ignoring distractors.  
114 The PBWM framework also accords with multiple lines of empirical evidence, ranging from neu-  
115 roimaging to manipulation studies, suggesting that the BG contributes to filtering (input gating)  
116 of WM (which improves effective capacity) (**McNab and Klingberg, 2008; Cools et al., 2007, 2010;  
117 Baier et al., 2010; Nyberg and Eriksson, 2016**) and selecting among items held in WM (output gat-  
118 ing; (**Chatham et al., 2014**)). Evidence also supports the PBWM prediction that striatal DA alters  
119 WM gating policies analogous to its impact on motor RL (**O’Reilly and Frank, 2006; Moustafa et al.,  
120 2008**); for review see (**Frank and Fossella, 2011**). Finally, these human studies are complemented  
121 by causal manipulations in rodent models implicating both striatum and thalamus as needed to  
122 support WM maintenance, gating and switching (**Rikhye et al., 2018; Nakajima et al., 2019; Wilhelm  
123 et al., 2023**). However, to date, these PBWM models have only been applied to WM tasks with dis-  
124 crete stimuli and thus have not addressed the tradeoff between precision and recall in VWM. Due  
125 to the discrete nature of the stimuli, accuracy is typically binary, and WM information could not be  
126 adaptively chunked. Further, previous PBWM studies only trained and tested within the allocated  
127 capacity of the model, limiting the application to common human situations in which set size of  
128 relevant items goes beyond WM capacity.

129 In sum, these two classes of neural WM models address complementary phenomena but their  
130 intersection has not been studied. In particular, how can our understanding of BG-PFC gating in-  
131 form the slots vs resources debate and the nature of WM capacity limits more generally? On the  
132 surface, PBWM is a slots model: it has multiple, isolated PFC “stripes” that can be independently  
133 updated and accessed for read-out. Note, however, that performance is improved in these models  
134 when they use distributed representations within the stripes (**Hazy et al., 2007; Kriete et al., 2013**),  
135 which can have resource constraints (**Frank and Claus, 2006**). We thus considered whether PBWM  
136 could acquire, via reinforcement learning, a gating strategy whereby it stores a “chunked” repre-  
137 sentation of multiple items within the same stripe, leaving room for other stripes to store other  
138 information and *in effect increasing the effective capacity without increasing the allocated capacity*.

139 Here, we sought to combine successful aspects of both models. We considered whether PBWM  
140 could be adapted to perform VWM tasks with continuous stimuli and whether it can learn a gating  
141 policy via RL that would support chunking to meet task demands. We include a ring attractor model

142 that allows for mergers of continuous-valued stimuli via overlapping representations (*Wei et al.*,  
143 *2012*; *Edin et al.*, *2009*; *Almeida et al.*, *2015*; *Nassar et al.*, *2018*). But rather than representing all  
144 input stimuli at once, the network evaluates a single sensory input in the context of stimuli already  
145 stored in one or more PFC stripes. The ring attractor can then merge or chunk the sensory input  
146 with its nearest neighbor in PFC. Importantly, the resulting chunks are not obligatorily stored in WM.  
147 Rather, the BG learns a gating policy so that it can potentially store the original (and more precise)  
148 sensory input, or it can replace a currently stored representation with the chunk – and adaptively  
149 alter its strategy as a function of task demands (set size). During recall, the network can gate out the  
150 corresponding (original or chunked) representation linked to the probe, and reproduce hallmarks  
151 of human performance in continuous report VWM tasks.

152 Notably, we find that this chunk-augmented PBWM network outperforms control models that  
153 lack chunking abilities across a range of task conditions. Chunk models outperform control net-  
154 works even when the control models are endowed with an allocated capacity that exceeds set size.  
155 This latter result stems from a credit assignment problem that arises when networks must learn  
156 to store and access multiple items in WM. Chunking instead allows for a common set of stripes  
157 to be repeatedly reused and reinforced, limiting the number of possible solutions explored. As  
158 such, this result lends insight into a normative rationale for why WM capacity is limited in the first  
159 place. Moreover, these performance advantages depend on a healthy balance of BG dopamine  
160 signaling needed to support adaptive gating policies that enhance effective capacity, providing a  
161 novel account for WM deficits resulting from aberrant BG DA signaling in patient populations such  
162 as Parkinson’s disease and schizophrenia (*Frank, 2005*; *Moustafa et al.*, *2008*; *Cools, 2006*; *Cools*  
163 *et al.*, *2010*; *Maia and Frank, 2017*). Finally, we show that, like humans, the model shows a recency  
164 bias, with increased accuracy for reporting items that had been presented more recently. This re-  
165 sults both from an increased propensity to update items that had been presented earlier and as a  
166 consequence of chunking of earlier items. This account is consistent with evidence in the human  
167 literature on the nature of recency effects *Oberauer et al. (2012)*, and inconsistent with alternative  
168 neural models of passive decay.

## 169 **Methods**

170 The model is implemented using an updated version of the Leabra framework (*O’Reilly et al.*, *2024*),  
171 written in the Go programming language (see <https://github.com/emer/emergent>). All of the com-  
172 putational models, and the code to perform the analysis, are available and will be published on our  
173 github account. We first outline the basic neuronal framework before elaborating on the PBWM  
174 implementation, modifications to the continuous report task, and chunking implementation, with  
175 most details in the Appendix.

176 Leabra uses point neurons with excitatory, inhibitory, and leak conductances contributing to  
177 an integrated membrane potential, which is then thresholded and transformed to produce a rate  
178 code output communicated to other units.

179 The membrane potential  $V_m$  is updated as a function of ionic conductances  $g$  with reversal (driv-  
180 ing) potentials  $E$  according to the following differential equation:

$$C_m \frac{dV_m}{dt} = g_e(t)\bar{g}_e(E_e - V_m) + g_l(t)\bar{g}_l(E_l - V_m) + g_i(t)\bar{g}_i(E_i - V_m), \quad (1)$$

181 where  $C_m$  is the membrane capacitance and determines the time constant with which the volt-  
182 age can change, and subscripts  $e$ ,  $l$  and  $i$  refer to excitatory, leak, and inhibitory channels respec-  
183 tively.

184 The excitatory net input/conductance  $g_e(t)$  is computed as the proportion of open excitatory  
185 channels as a function of sending activations times the weight values:

$$g_e(t) = \langle x_i * w_i \rangle = \frac{1}{n} \sum_i x_i w_i \quad (2)$$

186 Activation communicated to other cells ( $y_j$ ) is a thresholded ( $\Theta$ ) sigmoidal function of the mem-  
187 brane potential with gain parameter  $\gamma$ :

$$y = \frac{1}{\left(1 + \frac{1}{\gamma(g_e - g_e^\Theta)_+}\right)} \quad (3)$$

188 where  $g_e^\Theta$  is the level of excitatory input conductance that would put the equilibrium membrane  
189 potential right at the firing threshold  $\Theta$  and depends on the level of inhibition and leak.

$$g_e^\Theta = \frac{g_i(E_i - \Theta) + g_l(E_l - \Theta)}{\Theta - E_e} \quad (4)$$

190 Further details are in the appendix, but we elaborate the inhibition function below given its rele-  
191 vance for the chunking mechanism.

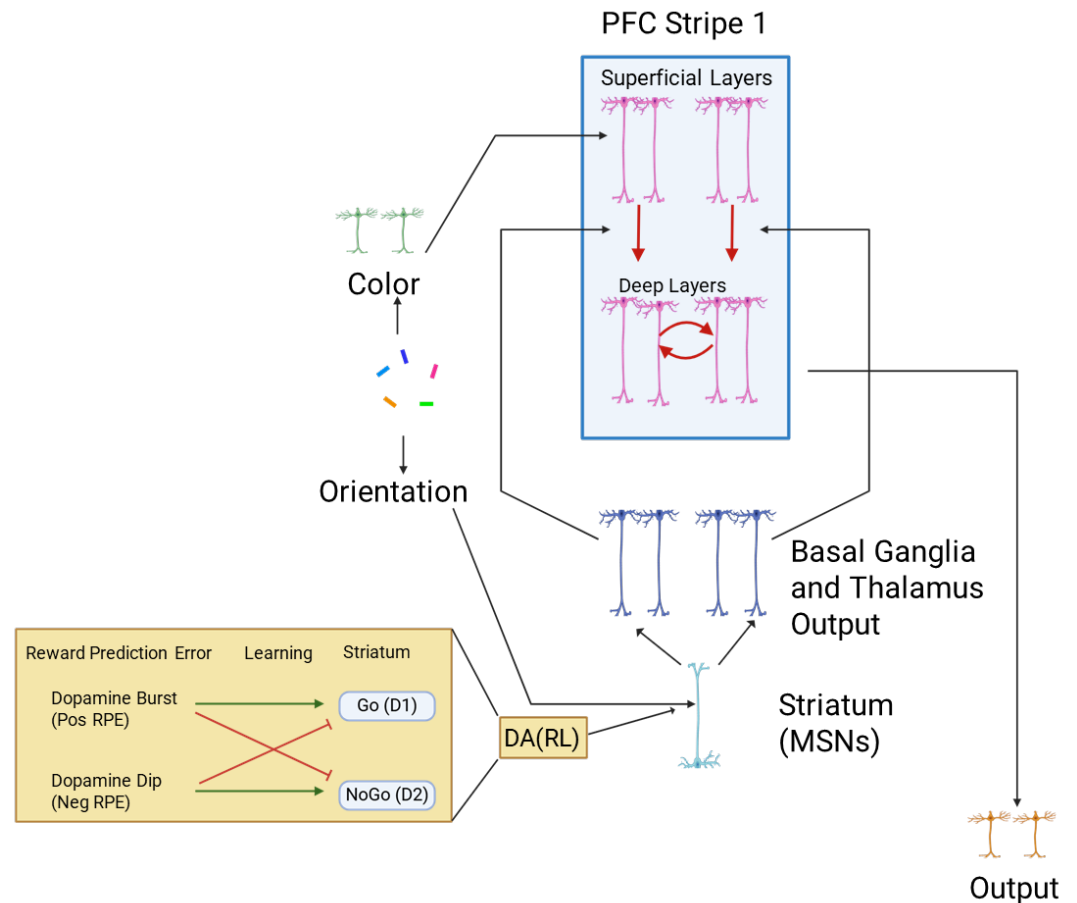
### 192 **Base PBWM model**

193 The PBWM model is built based on a large repertoire of research surrounding WM, gating, and  
194 RL and has been developed over a series of articles (see introduction). Here we focus on the high  
195 level description of its functionality and the unique additions to the current application, particularly  
196 the implementation of continuous rather than discrete representations throughout the network  
197 (input, PFC and response layers), and the chunking mechanism.

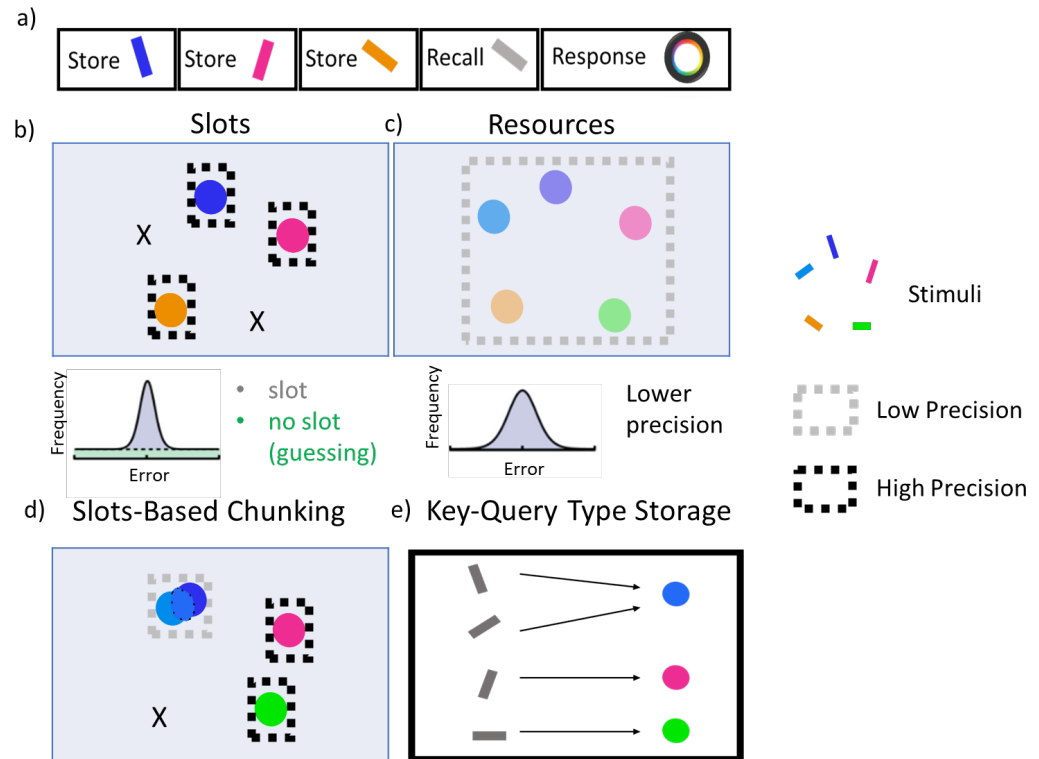
198 We modified PBWM to accommodate continuous representations such as those used in the de-  
199 layed report color wheel task (*Berg et al., 2012; Nassar et al., 2018*), see Figure 2a, a common task  
200 to assess precision and recall characteristics of VWM and which forms the main basis for our simu-  
201 lations. In this task, participants are presented with multiple colored bars on a visual display where  
202 each bar has two attributes: a color and orientation. After some delay, participants are shown only  
203 one of the orientations in gray, and the subject's task is to report the color that was associated with  
204 that orientation using a continuous color wheel. Previous PBWM applications used only discrete  
205 stimuli and did not address precision. To simulate continuous report tasks, we represented the  
206 color for each stimulus as randomly varying from 0 to  $2\pi$  using a population code, where each  
207 neuron in the layer maximally responds for a particular color, and the full representation for a  
208 given continuous input is represented as a Gaussian bump over an average of 10 neurons in a  
209 20 neuron layer. This representation accords with that seen in visual area V2, with hue neurons  
210 that are spatially organized according to color (*Xiao et al., 2003*). Each color was presented to the  
211 network together with a separate representation of its associated orientation (for simplicity we  
212 used discrete orientations, as the task is not to recall the precise orientation but only to use it to  
213 probe the color). The stimuli are presented sequentially to the mode. This serves 3 purposes: to  
214 simplify the binding problem, to mimic a sequential attentional mechanism, and to make contact  
215 with literature on serial WM tasks.<sup>1</sup>

216 These input layers project to the PFC maintenance layers, which contain isolated populations  
217 in discrete stripes, also coded as Gaussian bumps of activity. As in prior PBWM models, the PFC  
218 is divided into superficial and deep layers (*O'Reilly and Frank, 2006; Hazy et al., 2007; Frank and  
219 Badre, 2012; Hazy et al., 2021*). The superficial PFC layers for each stripe will always reflect the  
220 input stimuli transiently, as candidates to be considered for storage in WM. But for these stimuli to

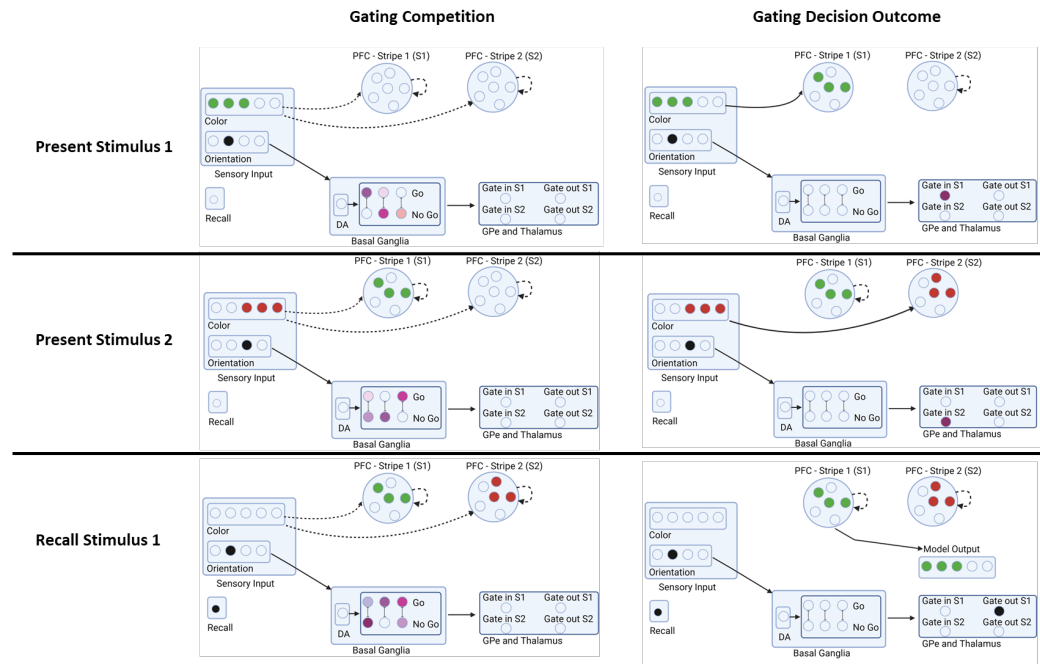
<sup>1</sup>While this sequential presentation of stimuli simplifies the binding problem, it also adds a further challenge as the network must have capacity to recall items that were presented several time steps/trials ago. To solve this problem, the model must learn an efficient and appropriate gating strategy. This also makes the model vulnerable to serial dependence in its chunking, consistent with that observed empirically, whereby WM reports are biased towards the last stimulus that was presented (*Bliss et al. (2017); Kiyonaga et al. (2017); Fischer and Whitney (2014)*)



**Figure 1. Base Model** Sensory inputs (reflecting visual cortical representations) project to PFC superficial layers, which transiently represent those inputs. Activity is maintained in PFC after stimulus offset (and into subsequent trials) only when gated. Red arrows indicate gating, supporting transfer of information from superficial layers to deep layers, triggered by striatal disinhibition of dorsomedial thalamocortical activity. Maintenance is represented by recursive red arrows in deep layer. Green insert shows how striatum D1 and D2 neural populations (which have opposite effects on gating) are modulated by dopaminergic reward prediction errors (RPEs). Over the course of learning, synaptic weights evolve to support effective gating strategies that increase rewards.



**Figure 2. Visual Working Memory Task.** a) The color wheel task is commonly used to study the nature of capacity limitations in VWM. During encoding, participants are presented multiple randomly generated oriented and colored bars. After a delay they are shown a recall probe trial in which one of the previously seen orientations is presented in gray. The participant responds by using a color wheel in an attempt to reproduce the color associated with that probe orientation. The number of store items are dictated by **set size**. b) Slots models suggest that WM capacity is limited by a fixed number of slots. When set size exceeds capacity, some items are stored in memory with high precision while the rest are forgotten, resulting in an error histogram that is a mixture of high precision memory (for items in a slot) and guessing (for items not in a slot). c) Resource models state that all items can be stored in a common pool, but as the number of items increase, the precision of each representation decreases, resulting in an error histogram with a large variance (but no guessing). Adapted from (*Ma et al., 2014*). d) A hybrid chunking model containing discrete slots, but with resource-like constraints within each slot. Here, the two bluish items are merged together within a slot, reducing their precision but freeing up other slots to represent pink and green items with high precision. The orange item is forgotten. The criterion for chunking can be adapted such that error histograms will look more like the slots theory or resource theory depending on task demands (WM load and chunkability of the stimulus array; (*Nassar et al., 2018*)). e) Storage in the PBWM-chunk model is like a key-query. The colors are stored as continuous representations in PFC and can be merged. The orientations are the queries used to probe where information should be stored and where to read it out from.



**Figure 3. Example Sequence of Network Gating Decisions.** In this example trial, the network is presented with stimulus 1 (color and orientation), stimulus 2, and is then asked to recall the color of stimulus 1 based on just its orientation. Each step is broken into a gating competition that involves the Basal Ganglia (striatal Go/NoGo units, Gpe/Gpi) and Thalamus units. The outcome of this internal competition determines the gating decision and the model output. When the first stimulus is presented, the relative activities determine if and where the stimulus is gated in (stripe 1 or stripe 2). The network gates stimulus 2 in a different stripe based on its orientation. During recall, the network uses a gating policy to output gate the stripe corresponding to the probed orientation. A reward is delivered to the network proportional to the accuracy in reporting the original color. A negative reward is delivered if the color is not sufficiently close (see Methods). Rewards translate into dopaminergic reward prediction error signals that serve to reinforce or punish recent gating operations. This schematic is illustrative; the actual network contains a broader population code and the PFC stripes are divided into input and output representations, each with deep and superficial layers (see Text).

221 be maintained robustly over delays (and over intervening other stimuli on subsequent time points),  
 222 they have to be gated into WM. Accordingly, each PFC stripe is modulated by a corresponding BG  
 223 module consisting of striatal "Go" and "NoGo" units which in turn, via direct and indirect pathways,  
 224 project to BG output / thalamic units. When there is relatively more Go than NoGo activity in a  
 225 given module, the corresponding Thalamic output unit is activated, inducing thalamocortical re-  
 226 verberation and activation of intrinsic ionic maintenance currents, thereby triggering robust main-  
 227 tenance of information in the deep PFC maintenance layers (*O'Reilly and Frank, 2006; Hazy et al.,*  
 228 *2007; Frank and Badre, 2012; Hazy et al., 2021*). Thus only the stripes that have been input gated  
 229 continue to maintain the most recent color representation in an attractor over time. Importantly,  
 230 gating in PBWM implements a form of "role-addressable" memory: the decisions about whether  
 231 and which stripe to gate colors into depends on its assigned role. In this case, the orientation probe  
 232 associated with the color determines where it should be gated. By receiving inputs from the ori-  
 233 entations, the BG can thus learn a gating policy whereby it consistently stores some orientations  
 234 into a particular PFC stripe, making it accessible for read out. Figure 3 shows a schematic example  
 235 in which, based on the orientation the network gates the first PFC stripe to store the green color,  
 236 but then stores the color of the second orientation to store the red color. Thus, the PBWM stripes  
 237 serve a variable binding function (*O'Reilly and Frank, 2006*) which can also be linked to key/query  
 238 coding (*Traylor et al., 2024; Swan and Wyble, 2014*): in this case the network can learn to use the  
 239 orientations to guide which stripe is accessed, and the population within the stripe represents the



240 content (color).

241 During a recall trial, the network is presented with only a single orientation probe. The model  
242 needs to correctly "output gate": select from multiple PFC maintenance stripes, so that only a sin-  
243 gle representation is activated in the corresponding PFC "output stripes", which in turn projects to  
244 the output layer (Figure 3 bottom row). If it correctly recalls the associated color (by activating the  
245 color population in the output layer) it receives a reward. Rewards were given in a continuously  
246 linear fashion based on the difference between output and the target color. The activity of the  
247 output neurons was decoded (using a weighted linear combination of neuron activities; (*Almeida*  
248 *et al., 2015*)) and the reward given was inversely proportional with the error. Correctly recalling a  
249 color thus requires not only having it stored in PFC, but reading out from the correct stripe that cor-  
250 responds to the probed item. This read-out operation involves a BG output gating function (*Hazy*  
251 *et al., 2007; Frank and Badre, 2012; Kriete et al., 2013*), facilitating transmission of information  
252 from a PFC maintenance stripe to the corresponding PFC output stripe. In this way, the model can  
253 read out from its several stored WM representations according to the current task demands (see  
254 (*Chatham et al., 2014*) for neuroimaging evidence of this BG output gating function). The input and  
255 output gating mechanism in PBWM performs a role-addressable gating function that can be linked  
256 to the key-query operations in Transformers ((*Traylor et al., 2024*)).

257 Note that successful performance in this and other WM tasks (*Hazy et al., 2007; Frank and*  
258 *Badre, 2012*) requires learning both the proper input gating strategies (which PFC stripes to gate  
259 information into, depending on the orientation, and which to continue to maintain during inter-  
260 vening items), and output gating strategies (which PFC stripes to read out from in response to a  
261 particular orientation probe). Such learning is acquired via reinforcement learning mediated by  
262 dopaminergic signals projecting to the striatum (represented by the bottom half of the model). At  
263 the time of recall, the network receives a dopamine (DA) burst or dip conveying a reward prediction  
264 error signal (RPE, by comparing the reward it receives with the reward it expected to receive using a  
265 Rescorla Wagner delta rule algorithm). These DA signals are used to modulate synaptic plasticity in  
266 the Go and NoGo units, reinforcing corticostriatal Go signals that drive adaptive input and output  
267 gating operations (following positive RPEs), and reinforcing NoGo units that suppress ineffective  
268 gating strategies (following negative RPEs). (See Appendix for information on parameter searching  
269 for other parameters of the model.) The model's previous gating decisions are flagged using an  
270 eligibility trace via synaptic tagging. Thus, Rewards modulate Go and NoGo synaptic weights not  
271 only based on their current activity levels but also based on these synaptic tags reflecting recent  
272 activity.

273 In addition to the gating strategies, which are learned via dopaminergic reinforcement learning  
274 as per above, the network also has to learn to produce the correct color in the output layer. On  
275 each Store and Ignore trial, regardless of whether stimuli are gated into PFC, the network has to  
276 report the color presented in the input. As such the connections from the input to output layers  
277 are plastic, and this mapping is learned via supervised learning (i.e., target colors are presented  
278 and the network learns from errors using the XCAL learning rule, see appendix). On recall trials,  
279 the network also has to learn to map the neurons that are output gated from PFC to reproduce the  
280 associated color in the output layer. As such, connections from PFC output stripes to the output  
281 layer are also plastic and learns according to the same rule. (Note that this is only useful if the  
282 corresponding stimuli have been gated in and out of PFC). All other weights are fixed (i.e., from  
283 striatum to GPiThal, and from PFC superficial to deep layers.

284 All networks are trained for 500 epochs and 100 trials per epoch. This was more than sufficient  
285 for learning to converge. We found that despite some fluctuations in training curves late in training,  
286 the gating strategy used by the model was largely unchanging at this point.

### 287 **Chunking / nearest neighbor implementation**

288 The base model incorporates purely feedforward visual input that can be stored in prefrontal cor-  
289 tex. But it is widely appreciated that there are also top-down projections from prefrontal cortex to

290 posterior sites (e.g., in parietal cortex). We posited that these top-down projections would allow  
291 networks to bias the sensory representation toward those stored in PFC, facilitating chunking. The  
292 degree of such bias should depend on the strength of those projections, and indeed experimen-  
293 tal evidence shows multiple abstractions of an item across parietal cortex and auxiliary sensory  
294 regions *Ito and Murray (2023)*.

295 We considered a minimal set up in which the network has access to two such representations,  
296 allowing it to represent both the raw sensory input (through connections directly to PFC) but also  
297 via a "chunking layer" representing posterior cortical association areas that receive excitatory input  
298 from both sensory regions and top-down information from deep PFC maintenance layers (Figure  
299 4a). This convergent excitatory connectivity ensures that if any of the currently maintained PFC  
300 representations is close enough to the current sensory input, the overlapping neural activations  
301 will be enhanced, thereby biasing the bump attractor in the chunk layer to be attracted toward  
302 the nearest PFC representation(s). Lateral inhibition ensures that only the most excited units will  
303 remain active. As such, the chunk layer will either represent the original incoming stimulus (if no  
304 PFC representation is sufficiently close to it) or a combined representation between the incoming  
305 and existing stimuli (Figure 4b).

306 More specifically, the excitatory conductance to the chunk layer comes from both the input  
307 layer and the PFC and can be summarized in the following equation:

$$g_e(t) = \frac{1}{n} \sum_i x_i w_i \quad (5)$$

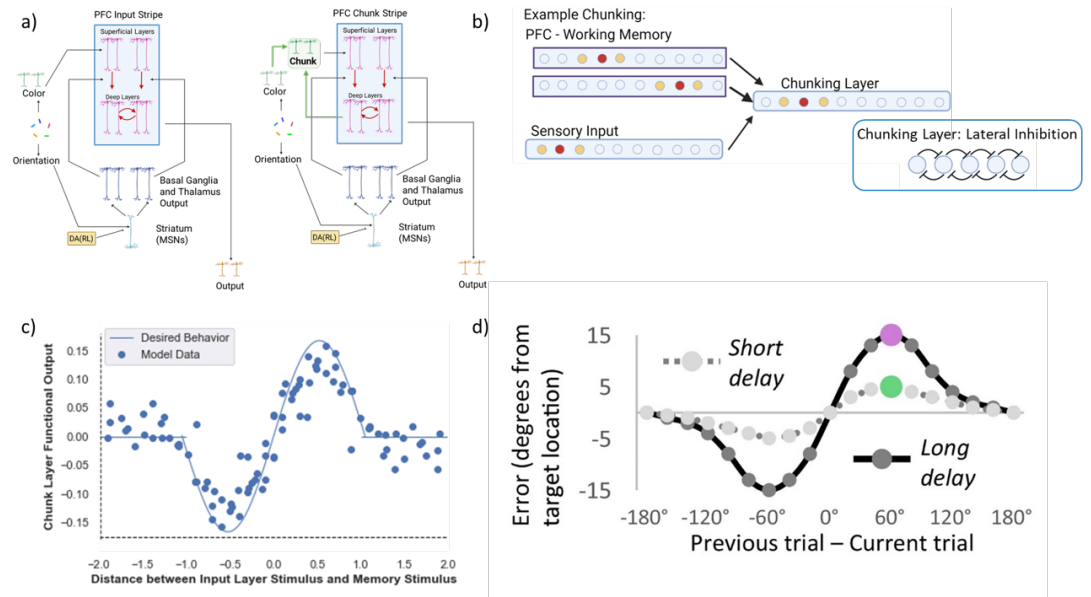
308 where  $g_e(t)$  is the excitatory conductance (input) to a layer,  $x_i$  is the activity of a particular sending  
309 neuron indexed by the subscript  $i$ ,  $w_i$  is the synaptic weight strength that connects sending neuron  
310  $i$  to the receiving neuron, and  $n$  is the total number of channels of that type (in this case, excitatory)  
311 across all synaptic inputs to the unit. Note that the relative strength of projections from input  
312 to PFC is scaled, with larger influences of input than PFC (to reflect e.g., number of synapses or  
313 proximity to the soma, using a relative weight scale ( $w_i$ )) (see appendix and *Computational Cognitive*  
314 *Neuroscience*, Chapter 2 (*O'Reilly et al., 2024*) for detailed information). This allows the chunking  
315 layer to preferentially reflect the input, subject to an attraction toward PFC representations that  
316 overlap with it (the nearest neighbor; Figure 4).

317 Lateral inhibition regulates activity in the chunking layer and, together with the strength of the  
318 top-down PFC projections, determines how close representations must be to bias the input toward  
319 the nearest PFC representations. (If the peak of the PFC bump attractor overlaps only with the  
320 tail of the sensory input bump, its influence will be largely suppressed due to inhibition). Lateral  
321 inhibition is implemented using feedforward (FF) and feedback (FB) inhibition (FFFB) which both  
322 alter the inhibitory conductance  $g_i(t)$ :

$$g_i(t) = G_i [ff(t) + fb(t)] \quad (6)$$

323 ,  
324 where feedforward inhibition ( $ff(t)$ ) is determined by the summed net input to the layer and  
325 feedback inhibition ( $fb(t)$ ) is determined by the firing rates of the neurons in that layer, and the  $G_i$   
326 gain parameter determines the overall sparsity of the layer (i.e., the relative influence of inhibition  
327 compared to excitation  $G_e$ ); see *O'Reilly et al. (2024)*.

328 Critically, a given PFC stripe receives projections from either the sensory input or the chunking  
329 layer. (The more general idea is that PFC stripes may have access to posterior cortical layers having  
330 varying levels of top-down projections and therefore chunking profiles). As such, PBWM could learn  
331 a gating strategy to store either the original sensory input into the corresponding PFC stripe or to  
332 store the chunked representation into the other stripe. In the latter case, it would replace the  
333 existing item stored in that PFC stripe with the new chunked representation, incorporating the  
334 novel input stimulus. These gating strategies are not hard-wired but are learned. For instance, the  
335 network could learn to use one stripe to store colors linked to particular orientations and use the



**Figure 4. PBWM Model and Chunking Layer Details.** a) Network diagram in the minimal case of two stripes: the first PFC stripe receives projections from the input layer ("PFC Input Stripe"); the second PFC stripe receives projections from the chunk layer ("PFC Chunk Stripe"). The network can be scaled to include more stripes of either type. We will refer to this model as the "chunk model". The control "no chunk" model has two stripes that are both of the type "PFC Input Stripe" (but it can use them to store separate inputs). b) Chunking schematic. A posterior ring attractor layer receives both bottom up sensory input and top-down input from the two PFC stripes (maintaining separate stimuli/representations). Overlap between the sensory input and the first PFC representation leads to convergent excitatory input in the chunking layer, resulting in a merged attractor. The impact of the more distant PFC representation is suppressed due to lateral inhibition. c) Chunking profile based on similarity. The x-axis shows the difference (in arbitrary units - comparable to radians) between the incoming stimulus and the nearest stimulus in PFC. The y-axis shows the deviation in the decoded chunk layer representation from the input stimulus. If the sensory input is close to a PFC representation, the chunk layer is attracted toward it. If the difference between the input stimulus and the nearest PFC representation is too large, the chunk layer largely mirrors the input (due to stronger input than PFC projections and lateral inhibition). This chunking profile closely matches that seen human memory representations, whereby memory reports are biased toward recent stimuli (top right inset, adapted from (Kiyonaga et al., 2017)).

336 other stripe for the rest of the orientations, allowing it to appropriately manage where to store and  
337 read out information when given the probe. In this case it would have precise memory for those  
338 representations that are stored and accessed, but it would have to guess if the probed item was  
339 not stored. At the other extreme, the model could learn to preferentially gate representations into  
340 and out of the chunk stripe but with less specificity. We will see later how the model gating policy  
341 depends on task demands (specifically set size) and evolves with learning.

342 For each experiment, at least 80 separate random seeds were run for each network, and results  
343 are averaged across them. (For select analyses requiring more observations, 160 separate random  
344 seeds were used). To test how set size affects learning and memory performance, the models  
345 were trained and tested with set sizes 2, 3, or 4. The set size determines the maximum number  
346 of stimuli that can be presented before recall trials. For example, set size 4 means that networks  
347 have to maintain up to 4 items before receiving a recall probe, and it may have to recall any of the  
348 preceding items.

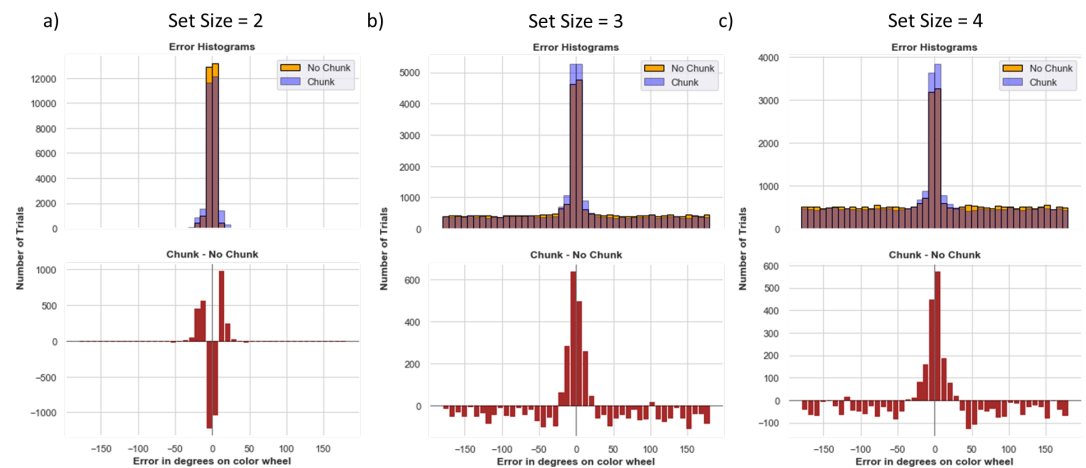
## 349 Results

350 We focus our simulations on variants of the color wheel task (see Methods). Briefly, networks  
351 were presented with continuous sensory inputs represented as coarse-coded Gaussian bumps of  
352 neural activity in conjunction with their associated discrete orientation. The number of stimuli to  
353 store ("set size") before a recall trial varied between 2 and 4 items.

354 During a recall "probe" trial, a single orientation was presented to the network, with no color  
355 (as in the empirical versions of this task). The model had to reproduce the associated color in the  
356 output layer, in the form of a continuous population-coded response. The error is then simply  
357 the difference between the decoded color from this population and the ground truth value that  
358 was presented during the earlier store trial for that orientation. "Correct" responses show as data  
359 points when errors are close to 0 degrees. If the model outputs an incorrect color, but that color  
360 corresponds to a different orientation, this is referred to as a binding or swap error (*Bays et al.,*  
361 *2009*). Finally, if the response is incorrect and nowhere near any of the colors for the stored orien-  
362 tations, it would be referred to as a guess. A guess could land anywhere along the axis of -180 to  
363 180 degrees and as such manifests as a uniform distribution across trials. If the model produced a  
364 non-response (nothing was gated out, e.g., if a stripe was empty), we randomly sampled from the  
365 uniform distribution ((*Almeida et al., 2015*) followed a similar process), mimicking random guess-  
366 ing. These errors (correct responses, guesses, binding errors) manifest themselves in the error  
367 distribution.

368 For proof of concept, we began with a minimal set-up in which all models were allocated 2 PFC  
369 maintenance stripes (we relax this assumption later to compare to models with larger allocated  
370 capacity). For all simulations, we compare performance (error distributions, gating strategies, in-  
371 fluence of dopamine manipulations) between networks with chunking ability (the "chunk model")  
372 against those with equivalent (or larger) number of stripes. We refer to the "allocated capacity" as  
373 the number of stripes given to the no-chunk model, because this is a hard limit on the maximum  
374 number of representations that can be stored and accessed. We refer to the "effective capacity"  
375 as the potentially larger number of items that can be accessed due to an efficient gating policy.  
376 Effective capacity can be improved if the network learns to consistently store (input gate) colors  
377 of distinct orientations in distinct stripes, and to appropriately read out from (output gate) the  
378 corresponding stripe at recall trial. It can also potentially be improved via chunking by increasing  
379 the number of representations that can be stored and accessed.<sup>2</sup> It is thus important to note  
380 that improving effective capacity requires an effective learning process to develop adaptive gating  
381 strategies, as the networks are not hard-coded to use any stripe for storing or accessing any repre-  
382 sentation. We will show how such learning depends on a healthy dynamic range of dopaminergic

<sup>2</sup>Note however, that effective capacity is not always larger than allocated capacity: without learning an effective input and output gating policy, a network's effective capacity will be less than its allocated capacity, for example if it overwrites information in the same stripe, if it accesses the incorrect stripe during recall, or if it doesn't gate a stimulus at all



**Figure 5. Model Recall Error Histograms.** The binned error in degrees is plotted on the x-axis, and the number of trials for that error bin on the y-axis. The blue and orange histograms show errors from all recall trials across all 80 random weight initializations from the chunk and no chunk models, each allocated with two stripes. The red histogram plots a bin-by-bin difference in errors between the models. a) For set size 2, there is very little difference between the models. The chunk model exhibits slightly higher rates of low errors neighboring zero (up to 30 degrees), due to small losses in precision resulting from some chunking (see text). b) Set size 3 is beyond the number of stripes allotted to the network. The chunk model has a larger density at zero and small errors, and less guessing (reduced density in the uniform distribution, see red lines). c) At set size 4, the chunking advantage is manifest by low errors and the improvement in less guessing becomes more pronounced (note y-axis scale - the reduction in guessing is actually reduced for set size 4 compared to 3).

**Figure 5—figure supplement 1. P(Recall) Across Set Size**

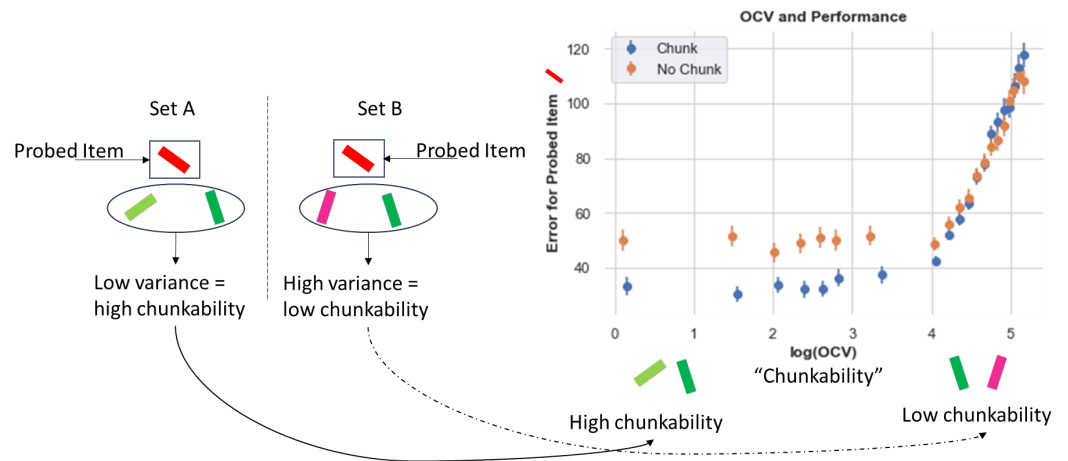
383 RL signals in the basal ganglia.

### 384 Error distributions across set sizes mirror those in human VWM and show chunking 385 advantages

386 Figure 5 shows a histogram of network errors (in degrees) during recall trials. The comparisons are  
387 made between chunk and no-chunk models as well as set sizes 2, 3, and 4; in these simulations we  
388 begin with a minimal setup in which both networks are endowed with two stripes. When set size is  
389 equal to the number of stripes (2), errors are small and centered around 0, with some variance due  
390 to precision. The overall performance is similar between models, but note that the chunk model  
391 shows somewhat more imprecise responses as indicated by some more small error trials. The  
392 ability to chunk results in a small cost which manifests as a decrease in precision when chunking  
393 occurred but did not need to be used.<sup>3</sup>

394 As set size increases beyond the allocated capacity, both models resort to guessing (random  
395 errors) on a subset of trials. This pattern is observable by the error histogram containing a mixture  
396 of errors centered around zero and a uniform distribution (see Figure 1), as commonly observed  
397 in humans (*Zhang and Luck, 2008*). Notably, the chunking model guesses less than the no chunk  
398 model and has a higher chance of precisely recalling the items (Fig 5). The difference between  
399 the chunk and no chunk model widens as the item limit continues to grow beyond the number of  
400 stripes. Comparing chunking and non-chunking models illustrates how the benefit of chunking is  
401 task-dependent. We next explored how chunking may improve the effective capacity of the net-  
402 work beyond the allocated number of stripes, and more specifically in which trials the advantage  
403 manifests, motivated by similar analysis in humans (*Nassar et al., 2018*). We subsequently explore

<sup>3</sup>Note that the chunk model is endowed with two stripes, and thus the only way for it to recall both items is to use both the input stripe and the chunk stripe. As a result, at least one item could be stored precisely in the input stripe, but if the other item is close enough to it, the PFC will store the less precise chunked representation in the chunk stripe, and memory reports will be biased. The network can nevertheless store two precise items if they are far enough apart such that the chunk layer is not biased by the other PFC representation (see Figure 4)



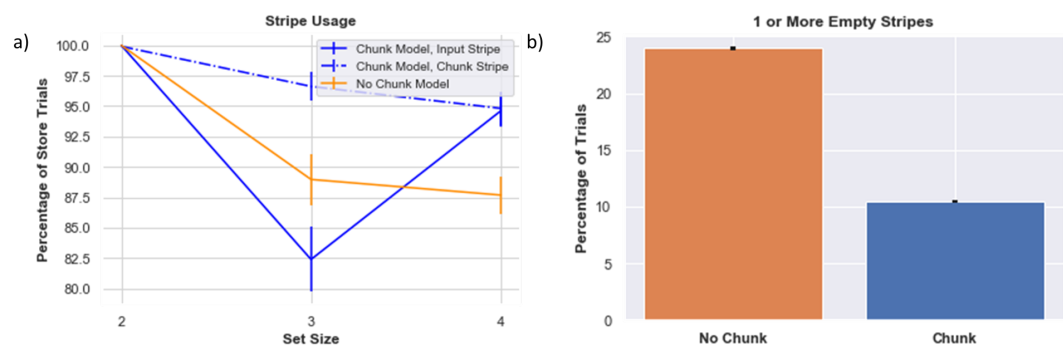
**Figure 6. Chunking improves recall for non-chunked items.** Left. Example array. Here we compare 2 sets, both containing a red item that will be later probed. In Set A, the other items (out of the probed cluster) are two shades of green and thus low variance (are similar to each other) and are therefore more likely to be chunked. In set B, the out of cluster variance (OCV) for the green and pink items is higher and these items are not likely to be chunked. Right. Chunking networks show consistent Recall advantages (lower errors) when OCV is low and hence the other items are chunkable. This difference disappears as OCV increases and overall errors rise. Errors plotted over all trials averaged over 80 networks in each case.

404 how these strategies are acquired via reinforcement learning.

### 405 **Chunking Frees Up Space for Other Items**

406 A key normative motivation for chunking is that doing so can save space to store other items,  
 407 thereby improving effective capacity (*Nassar et al., 2018*). In the model, when two items are chunked  
 408 into one stripe, the second stripe is free to hold another item. One key prediction is that  
 409 chunking should not only manifest in terms of loss of precision when probed with any of the chunked  
 410 items, but improved memory for the other items (Figure 2). Consider the situation in Figure 6  
 411 left, for a set size of 3. In both Set A and Set B, the red item is the probed item (the one that the  
 412 model is asked to recall). Set A contains other items in the set that are different shades of green  
 413 (low variance) and thus “chunkable”. If the network chunks them together, they will occupy one  
 414 stripe and the second stripe will be free for the red item, which is then more likely to be recalled at  
 415 high precision (as it is not part of the chunk). In Set B, the other items are pink and green and are  
 416 thus not chunkable (high variance). The network may store each of these into a separate stripe,  
 417 forcing it to randomly guess when probed with the red stimulus. (Alternatively, the red item could  
 418 be chunked with the pink item, in which case it will be recalled but with less precision than if it  
 419 stored only the original red item). To formalize and test this intuition, we quantified the “out of  
 420 cluster variance” (OCV) as the variance in degrees between the “other” items in the set (i.e., the  
 421 items that are not probed; see Appendix for details on OCV calculation). When this variance is low,  
 422 those items are more chunkable. We then assessed whether accuracy on recalling the probed item  
 423 (not part of this chunk) is improved as a proxy for chunking behavior.

424 Indeed, the data supports this prediction (Figure 6). When other items in a set are chunkable  
 425 (low OCV), the chunking network exhibits significantly smaller errors than the control network on  
 426 the probed item. After some OCV threshold, the chunking network no longer exhibits an advantage,  
 427 and both networks show increasing errors. (The increasing errors likely result from “swap errors”  
 428 (*Bays et al., 2009*), i.e., when the network reports one of the other stimuli in its memory instead of  
 429 the probed item - this results in larger errors as the entire set is more uniformly spaced and thus  
 430 not chunkable.) In sum, this analysis confirms that chunking does not merely result in reporting  
 431 imprecise responses for nearby items due to perceptual similarity, but that the network leverages  
 432 such similarity to its advantage so that it can save space for storing and recalling other items, as



**Figure 7. Stripe Usage** a) Stripe usage for the 1) chunk model, chunk -linked stripe 2) chunk model, input-linked stripe 3) no chunk model (average across both stripes), b) Proportion of trials when at least one stripe was empty. This analysis was done over 160 different model initializations.

433 also seen in humans (*Nassar et al., 2018*).

### 434 **Chunking leads to better resource management**

435 In addition to overall better performance, we hypothesized that chunk networks can manage mem-  
436 ory resources more efficiently than the no-chunk control models. We next compare how the mod-  
437 els use the allocated resources, focusing on store trials in which the maximum number of stimuli  
438 were presented (e.g., 4 stimuli for set size 4). We begin by analyzing the performance of networks  
439 after learning; below we explore how the chunk network learns to use the different gating strate-  
440 gies over the course of learning for different set sizes.

441 When the set size is 2, after learning, chunk and no-chunk models are equally able to utilize both  
442 stripes on 100% of the trials (Figure 7). The networks can properly gate both colors into distinct  
443 memory slots without overwriting or reusing the same stripe (in which case Fig 7 would have shown  
444 reduced use of one or the other stripe).

445 As the set size increases beyond allocated capacity, the gating management problem becomes  
446 much harder to solve via RL, and overall performance declines, particularly in the no-chunk (con-  
447 trol) model. Indeed, one might expect that as the number of items are increased, the model should  
448 always “max out” the number of stripes used, but in fact the opposite is the case. When the network  
449 attempts to gate information into both stripes, the no-chunk model will often receive negative re-  
450 ward prediction errors during recall trials when it will inevitably be forced to guess for a subset of  
451 the stimuli that is not in its allocated memory. As a result, input gating strategies will be punished,  
452 even if they were successfully employed and would have been useful had the other stimuli been  
453 probed. In turn, due to punishing gating policies that are generally useful, the stripe usage actually  
454 decreases, and the stripes are sometimes empty, akin to “giving up”. Conversely, if the network  
455 happened to be positively reinforced for gating a particular stripe, it might promiscuously gate  
456 that same stripe for other stimuli. This leads to overwriting information even though the other  
457 stripe is empty, and forcing the network to guess (or emit a non-response) when probed with a  
458 stimulus that was overwritten. In sum, the model exhibits a non-monotonic use of its resource, as  
459 its effective capacity actually declines relative to its allocated capacity. This result is reminiscent of  
460 experimental data in fMRI and EEG showing that PFC activity increases with increasing set size but  
461 then plummets when set size exceeds the participants capacity (e.g. (*Zhang et al., 2016*)), perhaps  
462 indicating a “giving up” strategy. We will explore the dopaminergic RL basis of such giving up in a  
463 section below.

464 In contrast, the chunk model is more effective at managing its resources when set size exceeds  
465 allocated capacity (Figure 7). Recall that the chunk model has access to one stripe with input from  
466 the chunked representation (“chunk stripe”) and one stripe with raw sensory input (“input stripe”).  
467 As set size increases, the network has more opportunities for chunking, and accordingly, relatively

468 more instances of reinforcing the gating operation linked to the chunk stripe. Interestingly, as set  
469 size just exceeds allocated capacity (here, for set size 3), the network decreases its use of the input  
470 stripe. This pattern arises because the network learns an advantage of chunking for set size 3 and  
471 thus sometimes does so more than it needs to (freeing up the input stripe), but also because of  
472 the cost of relying too much on the input stripe (as per the no-chunk model). Finally, as set size  
473 increases further (set size 4), the chunk network learns the benefit of storing some items in the  
474 held out input stripe, increasing effective capacity, while still effectively using the chunk stripe.

475 In sum, this analysis suggests that resource management and utilization is more consequential  
476 than the absolute number of stripes available. In previous network models of visual WM (*Wei*  
477 *et al., 2012; Nassar et al., 2018; Edin et al., 2009; Almeida et al., 2015*), responses were considered  
478 “correct” if the probed item was stored somewhere within the network; i.e., networks were not  
479 required to select among these stored representations in response to a probe, and they also did  
480 not have to decide whether or not to store an item in the first place. In contrast, the PBWM model  
481 focuses on the control of gating strategies into and out of WM, but requires RL to do so. Errors  
482 can result from reading out the wrong stripe (a swap error) or from an empty stripe (leading to  
483 guessing or non-responding). Chunking is an information compression mechanism that allows  
484 multiple stimuli to be mapped onto the same stripe. The chunk stripe has the advantage of being  
485 used repeatedly, giving the network has more opportunities to learn how and when to use that  
486 stripe.

### 487 **Chunking advantages remain even when comparing to networks with higher allo-** 488 **ated capacity**

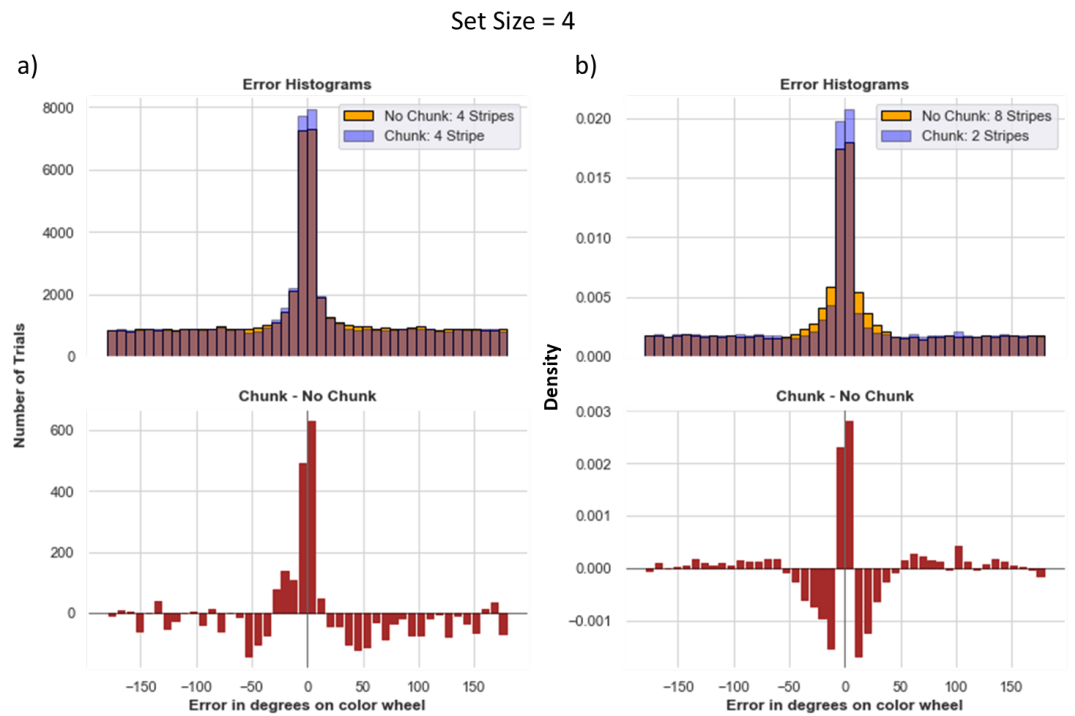
489 One might think that chunking advantages are limited to situations in which the allocated capacity  
490 is less than the set size. But when considering the challenges imposed in networks for which stor-  
491 age and access of items is not hard-wired, but must be learned, this is not a foregone conclusion.  
492 Indeed, above we found that the number of stimuli stored by no-chunk networks was even lower  
493 than their allocated capacity, due to RL challenges. We reasoned that such challenges could persist  
494 even when allocated capacity is increased to include or exceed the set size, due to credit assign-  
495 ment challenges in networks that are required to learn gating strategies into and out of multiple  
496 independent PFC stripes.

497 We begin with an analysis of networks performing the most difficult task (set size 4) but now  
498 allocated 4 stripes (for both chunk and non-chunk networks; in this case the chunk network has  
499 just one chunk stripe and 3 input stripes). The naive hypothesis states that no-chunk and chunk  
500 models should not differ in performance, or that chunk models might even show a deficit due to  
501 loss of precision when using the chunk stripe. Instead, based on earlier results and the above, we  
502 reasoned that chunk models might continue to show an advantage here, because frequent rein-  
503 forcement of gating into and out of the chunk stripe would reduce the credit assignment burden  
504 during learning that arises from learning to manage access of 4 different items into and out of WM.

505 Indeed, results supported these conclusions. Error distributions in chunk networks have more  
506 density at zero and low errors. The chunking model also guesses less (Figure 8).

507 This result largely persisted even in the more extreme case when allocating the control (no-  
508 chunk) model 8 stripes (twice as many as needed) and reverting the chunk model to using only 2  
509 stripes. While guessing is slightly reduced when having more stripes (due to more opportunities  
510 to store stimuli), the 8 stripe model does not increase the number of times it precisely recalls the  
511 correct item relative to the chunk model with only 2 stripes (Figure 8b). Upon further inspection,  
512 one can see that the 8 stripe model produced a larger density of moderate errors around the 0-30  
513 degree range. This result is curious because this model has no ability to chunk. To disentangle  
514 the source of these errors and to lend insight into the difficulty of WM gating strategy with high  
515 load and/or allocated capacity, we first removed trials in which the network did not give any output  
516 (these non-responses comprised roughly 12% and 24% of trials from the 8-stripe model and the  
517 2 stripe chunk model, respectively). Within the remaining trials, the 2-stripe chunk model has a





**Figure 8. Increasing Allocated Capacity  $\neq$  Better Performance: The importance of Resource Management** a) Chunk Model with 4 stripes vs. No chunk Model with 4 stripes in a task with set size 4. Even though the no-chunk networks has sufficient number of stripes to store each item with high precision, the corresponding chunk network still exhibits advantages, due to difficulties in credit assignment associated with managing four independent stripes. b) A more extreme comparison between the Chunk Model with 2 stripes vs. No chunk Model with 8 stripes. The chunk model guesses slightly more, but has more precise responses. The 8 stripe model has more density at small nonzero errors (see text for explanation). For both a and b the averages were computed over 160 models. For b, we display density rather than counts, because trials where either models gave no response were removed to better understand the small nonzero errors in the 8 stripe model (nonresponses add noise).

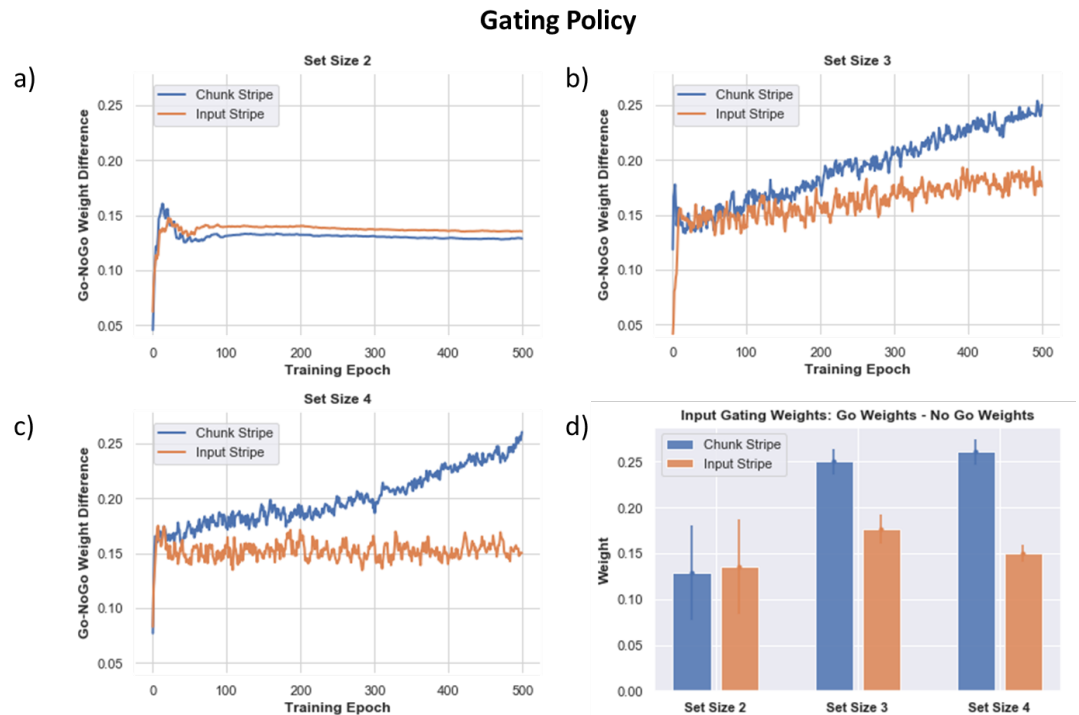
518 higher peak of precise responses (around 0 error) but still slightly more guessing than the 8 stripe  
519 no-chunk model, which continues to show a higher density at moderate errors (0-30 degrees). The  
520 reason for these moderate errors is that with 8 stripes, the network has a more difficult job. It  
521 needs to reinforce output gating strategies that properly read out from only the most appropriate  
522 stripe. Due to the curse of dimensionality (i.e., when the network outputs a response that is close  
523 enough to the target, it will get reinforced for any gating operations that preceded it, leading to  
524 spread of credit assignment to other stripes). Indeed, we found that the over-extended 8 stripe  
525 network frequently reads out from (output gates) multiple stripes in parallel (an average of 2.53  
526 stripes during recall), and thus even when the response does reflect the correct item it is also  
527 “contaminated” by reading out one of the other stripes, such that the averaged response results  
528 in a larger number of moderate errors. In contrast, the two stripe networks (across both no chunk  
529 and chunk) output gated an average of 1.08 stripes for each trial, appropriately reading out from  
530 a single PFC representation.

531 In sum, simply allocating a network with larger numbers of stripes does not yield the naïve  
532 advantages one might expect, at least when gating strategies need to be learned rather than hard-  
533 wired. In this case, the networks do use all the stripes available, but don't use them effectively. For  
534 example, qualitative observations revealed that a given network might gate one single stimulus into  
535 multiple stripes, and then proceed to overwrite many or all the same stripes with a new incoming  
536 stimulus – a strategy that is sometimes effective if it happens to get probed for the one of the  
537 items still in memory during recall. The large number of input and output gating operations to  
538 consider in tandem needed for adaptive behavior leads to a vexing credit assignment problem, as  
539 it becomes a challenge to know which of several gating operations or several PFC representations  
540 are responsible for task success / error, and networks fall into an unfortunate local minimum. This  
541 credit assignment problem is mitigated by chunking, allowing the network to reinforce the same  
542 input and output gating policy across multiple instances.

### 543 **Frontostriatal Chunking Gating Policy is Optimized via RL as a Function of Task** 544 **Demands**

545 The above results show that chunking networks confer advantages as set size grows, even com-  
546 pared to no-chunk networks that have a larger number of allocated stripes. Moreover, these ad-  
547 vantages come with little cost when set sizes are lower (e.g., 2). To explore how the network can  
548 adaptively learn to chunk as a function of task demands, we quantified the evolution over learning  
549 of each network's “gating policy”. Prior work has shown that PBWM develops a gating policy that  
550 predicts rapid improvement in task success when such policies mimic the task structure (e.g., for  
551 hierarchical tasks (*Frank and Badre, 2012*); see also (*Traylor et al., 2024*) who showed that modern  
552 transformer neural networks mimic a PBWM gating policy when challenged with WM tasks). Here  
553 we assessed whether networks could adaptively learn a gating policy that prioritizes gating into  
554 chunk vs input stripes depending on task demands.

555 In PBWM and related networks, the gating policy is dictated by learned synaptic weights into  
556 striatal GABAergic medium spiny neurons (MSNs). These MSNs are classically divided into D1 “Go”  
557 MSNs and D2 “NoGo” MSNs, with opponency between these populations determining which ac-  
558 tions are selected (i.e., those with the largest difference in Go vs NoGo activities; (*Frank, 2005*;  
559 *Jaskir and Frank, 2023*)). In the case of working memory, a network will be more likely to gate in-  
560 formation into a particular stripe if the synaptic weights are larger for the Go in comparison to the  
561 NoGo neurons. The relative weights control the disinhibition of that particular stripe. When the  
562 network performs well and it gets a reward prediction error, dopaminergic signals modify plasticity  
563 into the corresponding D1 MSNs, reinforcing the gating policy that drove the cortical update. Con-  
564 versely, errors associated with negative prediction errors lead to punishment of that gating policy  
565 by increasing synaptic weights into the D2 MSNs (*Frank, 2005*; *O'Reilly and Frank, 2006*). Below  
566 we confirm a key role for these dopaminergic signals in modulating adaptive performance. But  
567 first here we evaluated how they served to alter specific gating policies. We assessed PBWM gating



**Figure 9. Gating Policy (Go - NoGo Weights for Each PFC Stripe) Across Training** As the networks learn ( over 500 training epochs, averaged over 80 networks), the learned gating strategy differentiates between the input-linked (orange) or chunk-linked (blue) stripes. Positive values indicate the networks learn greater Go than NoGo weights for input gating stimuli into the corresponding stripe. a) Set size 2, the learned gating strategy shows a slight preference for the input stripe to be used (associated with increased precision), but the network also uses its chunk stripe to store the other stimulus (it is possible the chunk stripe stores a merged representation depending on the proximity of the stimuli) . b) As the set size increases to 3, the chunk stripe is increasingly preferred over training. c) This differentiation occurs earlier and more strongly for set size 4, where chunking has yet a larger advantage. d) Summary of Go - NoGo weights after training. A larger positive value shows a stronger preference for gating into that stripe. As set size increases, preference for gating into the chunk stripe increases. Relevant for training of all models: We can confirm that the network behavior has stabilized in learning even if the Go/NoGo weights continue to grow over time for the chunked layer (due to imperfect performance and reinforcement of the chunk gating strategy)

568 policies in terms of the differences in Go vs NoGo synaptic weights for each stripe, and how they  
569 evolved over time when networks were trained for each set size. Specifically, we computed the  
570 summed synaptic connection strengths from the Control Input Units representing Store Orienta-  
571 tions to the Go and NoGo input gating units in the PFC stripes corresponding to input or chunk:

$$GatingPolicy = \alpha_j = \left[ \frac{\sum Go - \sum NoGo}{\sum Go + \sum NoGo} \right]_+ \quad (7)$$

572 (Here  $[\ ]_+$  indicates that only the positive part is taken; when there is less Go than NoGo, the net  
573 input to the Thalamus is 0).

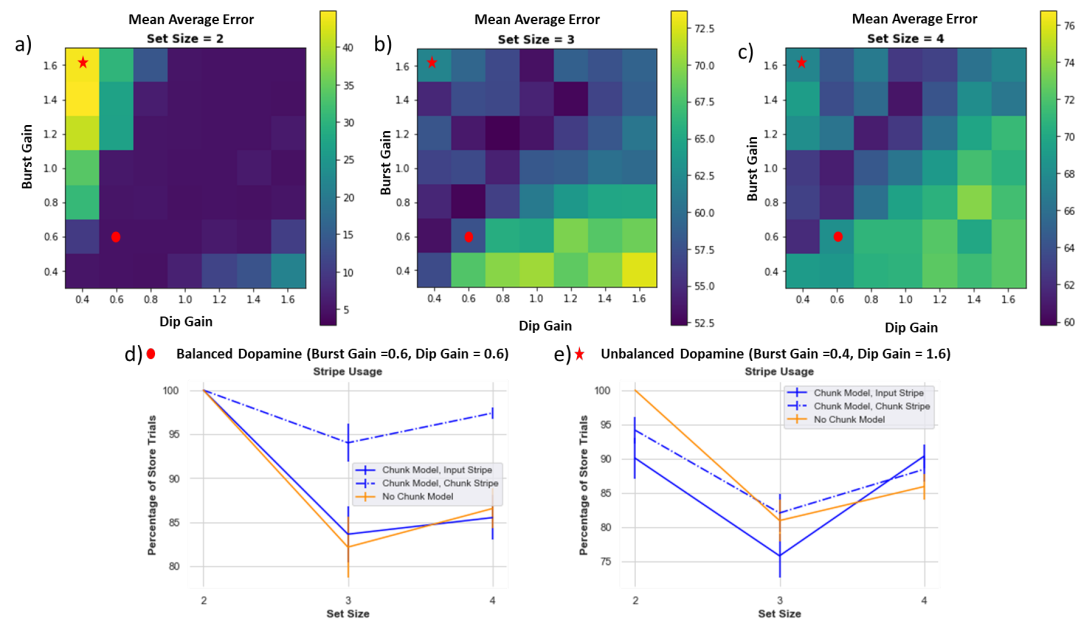
574 If the network learns that gating information into the chunk stripe is useful, it will evolve stronger  
575 Go vs NoGo weights for that particular gating action. But if instead it is more useful to gate the  
576 veridical input stimulus, it will develop a stronger gating policy for that stripe.

577 Figure 9 shows how such gating policies evolve over time. At set size 2 - where allocated ca-  
578 pacity (number of stripes) equals task demands, the gating policy slightly prefers to gate into the  
579 input stripe. This policy is sensible since the input stripe represents the original stimulus without  
580 any loss in precision, yielding lower errors. The network still learns a positive Go-NoGo gating pol-  
581 icy for the chunk stripe, because it can use that to represent the other stimulus. Notably, as set  
582 size increases, the chunk stripe increasingly becomes the preferred stripe for input gating over the  
583 course of learning. This adaptive change in gating policy allows the model to optimize a tradeoff  
584 between recall quantity and precision with increasing WM load, mediating the performance advan-  
585 tages described above. These results also accord with those observed in humans by (*Nassar et al.,*  
586 *2018*), whereby chunking propensities evolved adaptively as a function of reward history in their  
587 experiment, and also in their meta-analysis showing that chunking benefited performance more  
588 robustly in experiments with larger set sizes.

### 589 **Dopamine Balance is Critical to Learning Optimized Gating Strategies; Implications** 590 **for Patient Populations**

591 As noted above, learning gating strategies in PBWM is dependent on the basal ganglia and dopamin-  
592 ergic reinforcement system. Both chunk and no-chunk networks must learn whether to gate items  
593 into WM, which stripes to gate them into so that they can be later accessed (leaving maintained  
594 information in other stripes unperturbed), and during recall, which stripe should be gated out (de-  
595 pending on the probed orientation). To learn this, the network uses a simple RL "critic" which com-  
596 putes reward expectations and deviations thereof in the form of reward prediction errors (RPEs).  
597 Positive RPEs are signaled by dopamine bursts which reinforce activity-dependent plasticity in stri-  
598 atal Go neurons corresponding to recent gating operations (see Appendix and O'Reilly & Frank  
599 2006 for details). Conversely, when the model receives a reward that is worse than expected (i.e.,  
600 it reports an error), a dopamine dip (a decrease in phasic dopamine) will punish previous decisions.  
601 This negative RPE will punish the gating decisions by reinforcing corresponding NoGo neurons. To  
602 assess whether a healthy balance of such dopaminergic signals are needed for adaptive gating, we  
603 manipulated the gains of these dopaminergic bursts or dips to modulate their impact on Go and  
604 NoGo Learning. These investigations are relevant for assessing the basic mechanisms of the model  
605 but may also have implications for understanding well documented working memory impairments  
606 in patients with altered striatal dopaminergic signaling, such as Parkinson's disease, ADHD and  
607 schizophrenia (*Maia and Frank, 2017; Cools, 2006; Cools et al., 2007, 2008*).

608 Figure 10 shows how average absolute performance across 80 networks changes with DA burst  
609 and dip gain parameters. Overall, a healthy balance of relatively symmetrical DA bursts and dips  
610 is needed for optimized performance, but this effect also interacts with set size. The best perfor-  
611 mance for set size 2 (Figure 10a) is along the axis where burst and dip gain are symmetrical. As set  
612 size increases, the task becomes harder, and rewards are sparser due to more errors. In this case  
613 the best performance is on the axis where burst gain is somewhat greater than the dip gain; the  
614 model learns best when it can emphasize learning from sparse rewards.

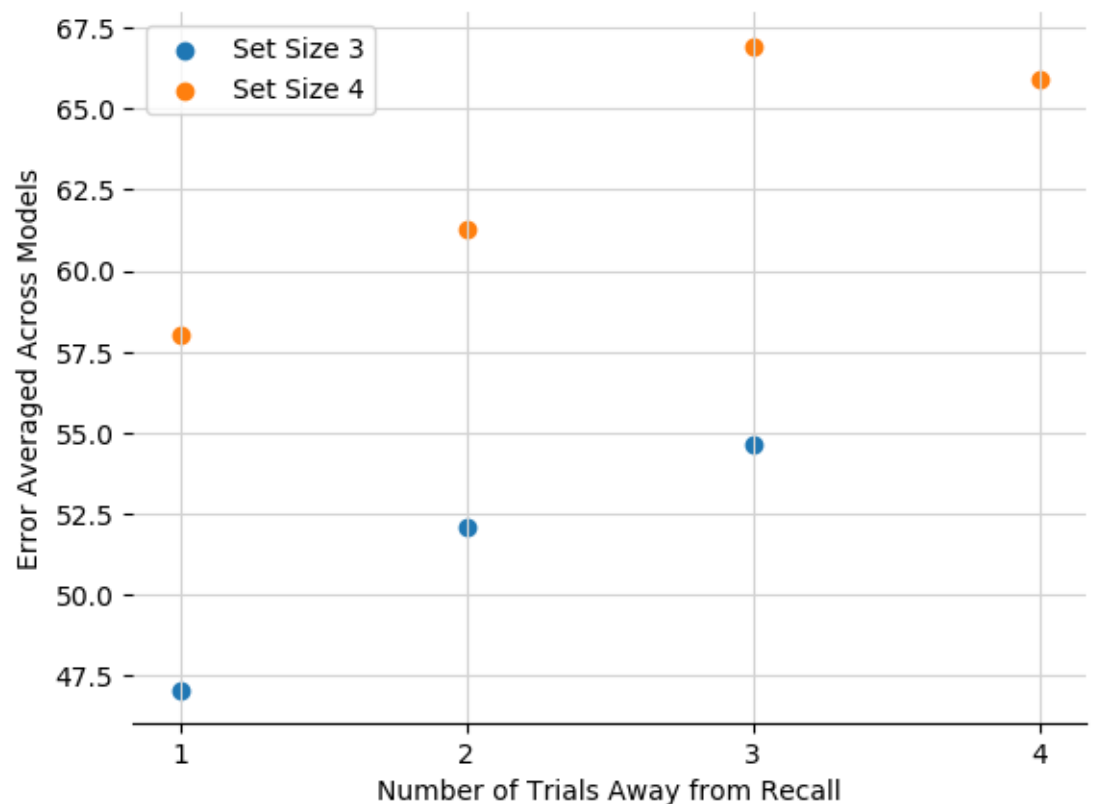


**Figure 10. Dynamic dopamine bursts and dips are needed for adaptive performance.** Each box is an average absolute error over 80 models. The color bar on the right indicates performance (note different scales on each plot), with darker colors (blue) representing better performance / lower absolute error. a) Set size 2: best performance is across the axis where burst and dip gain are symmetrical. b and c) Set size 3 and 4: best performance is where burst gain is slightly higher than dip gain. d) Example of balanced DA (burst gain = dip gain = 0.6), stripe usage. The chunk model manages to use the chunk stripe across all set sizes and both stripes in set size 2. The no-chunk model shows diminished storage of both stripes with increased set size due to greater propensity of DA dips. e) A regime of DA imbalance (larger DA dip than gain). The chunk model fails to robustly use both of its stripes, losing its advantage. The RL parameters interact with the ability for the chunk model to properly leverage chunking.

615 *Dopamine reinforcement signals can also lead to “giving up” and diminish effective capacity.* When  
616 set size increases, learning the proper gating strategies becomes difficult. The models may cor-  
617 rectly gate a few items in, but they may be incorrectly overwritten or incorrectly output gated at  
618 time of recall. Importantly, incorrect responses generate negative RPEs that punish preceding gat-  
619 ing actions, even if some of those actions were appropriate. A preponderance of negative RPEs  
620 can thus cause a network to “give up”, as observed in the no-chunk models when set size exceeds  
621 allocated capacity, leading to empty stripes (Figure 7). This mechanism is conceptually related to  
622 rodent findings in the motor domain, whereby lower DA levels can induce aberrant NoGo learning  
623 even for adaptive actions, causing progressive impairments (*Beeler et al., 2012*).

624 *Chunking can mitigate against giving up via shared credit assignment* The chunk model can combat  
625 against using the “giving-up” strategy: when items are chunked, the chunk stripe is used more  
626 frequently and therefore has a greater chance of receiving the positive reinforcement, and thus  
627 benefits from shared credit assignment. The model still has to learn when to use the chunk vs.  
628 input stripes, but chunking serves as an aid to the reinforcement learning process. Therefore,  
629 the chunk model is also more robust compared to the no-chunk model across various parameter  
630 ranges of dopamine. However, the chunk model still fails if the DA dip value is sufficiently larger  
631 than the DA burst, for similar “giving up” reasons (Figure 10 c,e).

### 632 **Network recapitulates human sequential effects in working memory.**



**Figure 11. Network Captures Recency Effects.** Average error on recall trials as a function of the distance in trials between presentation of the relevant stimulus and recall.

633 Finally, we also tested whether the model can reproduce well documented sequential effects  
634 in human working memory studies. Various findings indicate that in WM experiments, humans  
635 show higher accuracy for items most recently encountered. Moreover, *Oberauer et al.* shows  
636 how recency effects in humans stem not just from passive decay but specifically from intervening

637 distractors. Our network reproduces these effects, with average error monotonically increasing  
638 with number of intervening trials since the relevant stimulus was encountered, in both set sizes  
639 3 and 4. This results from the simple principle whereby with more intervening trials, the gating  
640 model is more likely to have updated a corresponding stripe, either replacing it altogether (leading  
641 to increased probability of forgetting) or chunking with previous items, thereby increasing average  
642 error. The errors are overall smaller in set size 3 because the network is less likely to overwrite an  
643 item altogether (it can chunk and still recall one of the items perfectly). Finally, the recency effects  
644 asymptotes in set size 4 due to increased opportunities for chunking (indeed, no-chunk networks  
645 continued to show increasing errors 4 trials back; not shown).

## 646 Discussion

647 Our work synthesizes and reconciles theories of visual WM across multiple levels of abstraction.  
648 At the algorithmic level, there has been extensive debate regarding the nature of WM capacity  
649 limitations, with experimental results alternately supporting slots or resources theories (*Bays et al.,*  
650 *2009; Berg et al., 2012; Wei et al., 2012; Swan and Wyble, 2014*). At the mechanistic level, several  
651 studies across species and methods suggest that circuits linking frontal cortex with basal ganglia  
652 and thalamus support WM input and output gating (*McNab and Klingberg, 2008; Cools et al., 2007,*  
653 *2010; Baier et al., 2010; Nyberg and Eriksson, 2016; Chatham et al., 2014; Wilhelm et al., 2023;*  
654 *Rikhye et al., 2018; Nakajima et al., 2019*). These data accord with the PBWM and related neural  
655 network models of PFC-BG gating processes (*Frank et al., 2001; O'Reilly and Frank, 2006; Hazy*  
656 *et al., 2007; Krueger and Dayan, 2009; Stocco et al., 2010; Badre and Frank, 2012; Calderon et al.,*  
657 *2022*). To date, however, these lines of literature have for the most part not intersected. Here we  
658 show that when augmented with continuous population code values and a chunking layer, PBWM  
659 comprises a hybrid between slots and resources with resource-like constraints within individual  
660 stripes. Moreover, through reinforcement learning, adaptive gating policies can adjust the degree  
661 to which behavior mimics primarily slots-like or resource-like as a function task demands. As such,  
662 this model accounts for human findings supporting chunking of related items in WM, that such  
663 chunking evolves with reward feedback, and is predictive of better performance with increasing  
664 task demands across multiple datasets (*Nassar et al., 2018*).

665 On its surface, PBWM is ostensibly a slot-like model, with distinct PFC clusters ("stripes") corre-  
666 sponding to slots that can be independently gated, giving rise to useful computational properties  
667 such as variable binding, indirection, compositionality and hierarchical generalization (*O'Reilly and*  
668 *Frank, 2006; Kriete et al., 2013; Collins and Frank, 2013; Calderon et al., 2022*). The need for gat-  
669 ing also accords with data suggesting that effective WM capacity is related to management of WM  
670 content, ie. one's ability to filter out distractors so as to prioritize task-relevant information (*Vo-*  
671 *gel et al., 2005; Astle et al., 2014; Feldmann-Wüstefeld and Vogel, 2019*), and that such abilities  
672 rely on basal ganglia output and function (*McNab and Klingberg, 2008; Baier et al., 2010*). How-  
673 ever, previous applications of PBWM and related networks have not confronted tasks requiring  
674 storage of continuous valued stimuli. In this work we augmented PBWM with a ring attractor layer  
675 that resembles that which has previously been applied to such continuous report tasks, supporting  
676 population level coding and mergers of nearby attractors (*Wei et al., 2012; Edin et al., 2009; Nassar*  
677 *et al., 2018*). However, in our network, this layer receives bottom-up input from both the sensory  
678 input layer and top-down input from all the PFC stripes, thereby allowing the network to combine  
679 sensory information with the nearest neighbor in memory. Moreover, WM chunking in our model  
680 is not obligatory, as the network can learn a gating policy that prioritizes raw sensory inputs (in  
681 which case it can represent a given item precisely) or to either replace a currently stored PFC item  
682 with the chunked version. As such, the PBWM-chunk model can learn to store more items than  
683 the allocated slots capacity by combining representations while incurring the cost of lost precision  
684 on the chunked items, giving rise to resource-like behavior. Given this learned policy, the network  
685 still may encounter trials where chunking is not possible and all stripes are occupied, leading to  
686 guessing and slots-like behavior. Depending on the learned gating policy and the task, the errors

687 look more "slots-like" or "resource-like".

688 As such, our model addresses key limitations in previous neural models in this domain, in which  
689 chunking was obligatory fashion due to perceptual overlap and could not be optimized (*Wei et al.,*  
690 *2012; Nassar et al., 2018*). Instead, PBWM adapts whether or not to chunk as a function of task  
691 demands and reward history (Figure 9), similar to empirical data (*Nassar et al., 2018*). Further,  
692 PBWM can also report only the color of the probed item, unlike previous neural models which  
693 were considered accurate as long as the probed color was one of the various populations still  
694 active (*Wei et al., 2012; Nassar et al., 2018*).

695 Critically, to perform adequately, PBWM requires learning appropriate input and output gating  
696 policies which are not hard-wired, and indeed involves solving a difficult credit assignment problem  
697 (*O'Reilly and Frank, 2006*). At the input gating level, the network must learn whether to gate the  
698 chunked or the raw sensory stimulus (via updating of the chunk vs input stripe). Simultaneously  
699 it must also learn which stripe to output gate in response to a given probe, which requires coor-  
700 dinating its input and output gating strategies so that they align. The credit assignment problem,  
701 understanding which input gating decisions in combination with output gating decisions lead to  
702 reward, is difficult. To understand the difficulty, we can look at an example case where the model  
703 input gates into stripe 1. However, during read out, since it has not learned the proper binding  
704 yet, it gates out from stripe 2, leading to an error and a dopamine dip. In this case, due to an  
705 improper output gating decision, both input gating decisions and output gating decisions will be  
706 punished. Eventual successful performance requires exploration of alternate gating strategies and  
707 reinforcement of those that are effective.

708 How can chunking help? First it is important to note that the above problem becomes even  
709 more difficult as the number of stripes increases – even if it matches or exceeds the set size (as  
710 shown in Figure 8). For example, random exploration and guessing will lead to the correct response  
711 (an item being gated into a stripe AND read out from the correct stripe) 50% of the time with 2  
712 stripes and 33% if the model has 3 stripes. The general form is:  $\frac{(n-1)^{N-1}}{n^N}$  where n is the number of  
713 stripes and where N is the set size. For a quick intuition, we assume that the first item is gated into  
714 any one of the stripes. We then multiply two probabilities: 1) the probability that the second item  
715 is gated anywhere *except* where the first item was stored - which is  $n - 1/n$  for 1 additional item.  
716 This probability is multiplied as many times based on the size size minus 1 since the first item is  
717 already stored (the power is  $N - 1$ ) 2) The probability that the first item is gated out correctly, which  
718 is  $1/n$ . The probability of this correct guess goes down as the number of stripes increases. The  
719 difficulty also increases with set size because the network must learn where to input and output  
720 gate for each item, and it is also possible for it to overwrite information by updating a stripe. As  
721 such, using a smaller number of stripes but allowing for chunking provides a lossy compression  
722 strategy that can mitigate this problem and render credit assignment easier, despite the loss of  
723 precision. Rather than overwriting information, the network can learn to use the chunk-linked PFC  
724 stripe if it is predictive of task success (minimal cost in the reward function for small errors), and  
725 moreover, when chunking is "good enough" the network can leverage repeated reinforcement of  
726 the same gating policy to store and read out from the chunked-link PFC stripe, thereby improving  
727 credit assignment.

728 As such, our simulations provide a provocative if somewhat speculative understanding of the  
729 nature of WM capacity limitations. It is unlikely that such limitations result from limits in the num-  
730 ber of neurons or stripes available in prefrontal cortex, given that discrete estimates on capacity  
731 limitations range in the order of 3-4 items whereas the number of stripes (or equivalent clusters  
732 of PFC populations) is orders of magnitude larger (*Frank et al., 2001*). Our simulations show that a  
733 limiting step is properly utilizing and managing resources to optimize performance (Figure 7), and  
734 that it might actually be more effective to limit the number of representations used. Increasing  
735 model capacity to 4 and 8 stripes and the resulting comparisons show that the limitation in the  
736 model is not simply about number of slots but the complexity of learning. Using WM requires mul-  
737 tiple input and output gating decisions and strategies in tandem with solving and learning a task



738 - this would become trivial with a homunculous dictating what information to store and where to  
739 store it. In biology, the PFC has to *learn* these gating strategies: it is not hardwired. This set of  
740 experiments helps explain various other WM findings which suggest that effective WM capacity  
741 is not just about "capacity" but rather is also about the ability to filter out irrelevant information,  
742 the importance of the task (reward), and experience with the task (experts vs. novice) (*McNab  
743 and Klingberg, 2008; Astle et al., 2014; Nassar et al., 2018; Feldmann-Wüstefeld and Vogel, 2019;  
744 Nakajima et al., 2019*). It also accords with our findings that when exceeding capacity, networks  
745 often "gave up" in the sense that they had more trials in which at least one stripe was empty, due  
746 to the influence of negative prediction errors punishing gating policies. As such, we showed that  
747 networks require a larger dopamine burst than dip to succeed with increasing task demands. This  
748 finding also accords with related data in rodents and our network model in the motor domain,  
749 whereby dopamine depletion can cause a progressive 'unlearning' of adaptive strategies (i.e., "giv-  
750 ing up") via aberrant NoGo learning (*Beeler et al., 2012*). This learned Parkinsonism was shown to  
751 be related to plasticity in D2 medium spiny neurons (*Beeler et al., 2012*), and this mechanism was  
752 recently confirmed to depend on the indirect pathway (*Cheung et al., 2023*).

753 More generally, our simulations revealed an important role for RL in shaping gating policies as  
754 a function of task demands, mimicking normative analysis showing that optimal chunking criterion  
755 changes with set size (*Nassar et al., 2018*). In the network, dopamine is a critical component of RL  
756 by adjusting synaptic weights into striatal modules that support input and output gating. The need  
757 for a healthy balanced dynamic range of DA signals for adaptive performance provides a potential  
758 window into a mechanism that can explain deficits in patient populations with altered striatal DA  
759 signaling. Whereas much of the literature in patients with schizophrenia and ADHD focuses on  
760 limitations in WM capacity, our simulations suggest an alternative whereby altered DA signaling in  
761 these populations (*Maia and Frank, 2017*) could influence chunking and efficient use of resources.  
762 Our finding that adaptively learned gating policies are important for controlling when and whether  
763 to chunk may have implications for recent accounts of patient populations with repetitive negative  
764 thinking, which are proposed to arise from failures inherent to learning adaptive mental gating  
765 (*Hitchcock and Frank, 2024*). Future work in these patient populations could aim to study these  
766 nuances for better understanding of their cognitive deficits.

## 767 **Limitations and Future Directions**

768 There are several limitations to this work. For simplicity, we restricted our simulations to a chunk-  
769 ing network with just one chunk-linked PFC stripe and one or more input stripes. In this case, the  
770 determining factor for whether stimuli are merged in the chunking layer depends on how close  
771 they are in color, lateral inhibition in the chunking layer, and the relative strength of top-down PFC  
772 projections to the chunk layer. These parameters were fixed in our simulations and were not for-  
773 mally optimized. A more general model could include multiple chunking layers with a reservoir of  
774 effective chunking thresholds (e.g. with varying degrees of top-down influence and lateral inhibi-  
775 tion). Depending on the task, the model could learn to chunk more liberally (larger set size - larger  
776 threshold) or more restrictively (smaller set size - smaller threshold), by adapting gating policies to  
777 rely on PFC stripes linked to these finer or coarser representations. Alternatively, it is possible that  
778 a network could learn to adapt these hyperparameters directly within the chunking layer. Further,  
779 through development the brain learns the environmental statistics and could learn those thresh-  
780 old parameters on a developmental time scale and could be fine-tuned on a task by task basis. Our  
781 objective was to explore how far one can get by optimizing only the gating policy via biologically  
782 plausible RL rules explored widely in basal ganglia.

783 Because of its wide application in the literature, we considered tasks in which stimuli can be  
784 chunked along a single scalar dimension (color or orientation, both of which have shown evidence  
785 for chunking *Nassar et al. (2018)*). Future work should explore to what degree these principle  
786 could generalize to more complex stimuli where chunking could occur across other more abstract  
787 dimensions, depending on the task demands (*Kiyonaga et al., 2017*). This model has the potential

788 to be scaled up and here we show the core principles for how the chunking gating strategy can be  
789 learned via RL.

790 One key difference is how the task is presented to the model and to humans. Humans are  
791 given clear verbal instructions and are able to perform the color wheel task with little to no prac-  
792 tice. However, the model does not receive verbal communication and must learn the task from  
793 scratch - random weights. It has no prior experience with how to process the stimuli or how to  
794 maintain any stimuli. Humans learn this through development and establish a general gating pol-  
795 icy. In a everyday setting, while individuals are not re-learning connections and gating policies to  
796 fit individual tasks, they are "fine-tuning" how they manipulate the information in WM and how to  
797 apply their learned policies to adapt to the current task. Experimental results show how reward or  
798 task relevance are factors that can tweak gating policies ((*O'Reilly and Frank, 2006; Nassar et al.,*  
799 *2018*)).

#### 800 **Acknowledgments**

801 Aneri Soni was supported by NIMH training grant T32MH115895 (PI's: Frank, Badre, Moore).  
802 The project was supported by ONR MURI Award N00014-23-1-2792 and NIMH R01 MH084840-08A1.  
803 Computing hardware was supported by NIH Office of the Director grant S10OD025181.

#### 804 **References**

- 805 **Almeida R**, Barbosa J, Compte A. Neural circuit basis of visuo-spatial working memory precision: A computa-  
806 tional and behavioral study. *Journal of Neurophysiology*. 2015; 114. doi: [10.1152/jn.00362.2015](https://doi.org/10.1152/jn.00362.2015).
- 807 **Astle DE**, Harvey H, Stokes M, Mohseni H, Nobre AC, Scerif G. Distinct neural mechanisms of individual and de-  
808 velopmental differences in VSTM capacity. *Developmental psychobiology*. 2014; 56. doi: [10.1002/dev.21126](https://doi.org/10.1002/dev.21126).
- 809 **Badre D**, Frank MJ. Mechanisms of hierarchical reinforcement learning in cortico-striatal circuits 2: Evidence  
810 from fMRI. *Cerebral Cortex*. 2012; 22. doi: [10.1093/cercor/bhr117](https://doi.org/10.1093/cercor/bhr117).
- 811 **Baier B**, Karnath HO, Dieterich M, Birklein F, Heinze C, Müller NG. Keeping memory clear and stable - The  
812 contribution of human basal ganglia and prefrontal cortex to working memory. *Journal of Neuroscience*.  
813 2010; 30. doi: [10.1523/JNEUROSCI.1513-10.2010](https://doi.org/10.1523/JNEUROSCI.1513-10.2010).
- 814 **Bays PM**, Catalao RFG, Husain M. The precision of visual working memory is set by allocation of a shared  
815 resource. *Journal of Vision*. 2009; 9. doi: [10.1167/9.10.7](https://doi.org/10.1167/9.10.7).
- 816 **Beeler JA**, Frank MJ, McDaid J, Alexander E, Turkson S, Bernandez MS, McGehee DS, Zhuang X. A Role for  
817 Dopamine-Mediated Learning in the Pathophysiology and Treatment of Parkinson's Disease. *Cell Reports*.  
818 2012; 2. doi: [10.1016/j.celrep.2012.11.014](https://doi.org/10.1016/j.celrep.2012.11.014).
- 819 **van den Berg R**, Awh E, Ma WJ. Factorial comparison of working memory models. *Psychological Review*. 2014;  
820 121. doi: [10.1037/a0035234](https://doi.org/10.1037/a0035234).
- 821 **Berg RVD**, Shin H, Chou WC, George R, Ma WJ. Variability in encoding precision accounts for visual short-term  
822 memory limitations. *Proceedings of the National Academy of Sciences of the United States of America*. 2012;  
823 109. doi: [10.1073/pnas.1117465109](https://doi.org/10.1073/pnas.1117465109).
- 824 **Bliss DP**, Sun JJ, D'Esposito M. Serial dependence is absent at the time of perception but increases in visual  
825 working memory. *Scientific Reports*. 2017; 7. doi: [10.1038/s41598-017-15199-7](https://doi.org/10.1038/s41598-017-15199-7).
- 826 **Brady TF**, Konkle T, Alvarez GA. A review of visual memory capacity: Beyond individual items and toward  
827 structured representations. *Journal of Vision*. 2011 5; 11:4-4. doi: [10.1167/11.5.4](https://doi.org/10.1167/11.5.4).
- 828 **Brady TF**, Alvarez GA. Contextual effects in visual working memory reveal hierarchically structured memory  
829 representations. *Journal of Vision*. 2015; 15. doi: [10.1167/15.15.6](https://doi.org/10.1167/15.15.6).
- 830 **Calderon CB**, Verguts T, Frank MJ. Thunderstruck: The ACDC model of flexible sequences and rhythms in  
831 recurrent neural circuits. *PLoS Computational Biology*. 2022; 18. doi: [10.1371/journal.pcbi.1009854](https://doi.org/10.1371/journal.pcbi.1009854).
- 832 **Chatham CH**, Frank MJ, Badre D. Corticostriatal output gating during selection from working memory. *Neuron*.  
833 2014; 81. doi: [10.1016/j.neuron.2014.01.002](https://doi.org/10.1016/j.neuron.2014.01.002).

- 834 **Cheung THC**, Ding Y, Zhuang X, Kang UJ. Learning critically drives parkinsonian motor deficits through imbal-  
835 anced striatal pathway recruitment. *Proceedings of the National Academy of Sciences of the United States*  
836 *of America*. 2023; 120. doi: [10.1073/pnas.2213093120](https://doi.org/10.1073/pnas.2213093120).
- 837 **Collins AGE**, Frank MJ. Cognitive control over learning: Creating, clustering, and generalizing task-set structure.  
838 *Psychological Review*. 2013; 120. doi: [10.1037/a0030852](https://doi.org/10.1037/a0030852).
- 839 **Cools R**, Miyakawa A, Sheridan M, D'Esposito M. Enhanced frontal function in Parkinson's disease. *Brain*. 2010;  
840 133. doi: [10.1093/brain/awp301](https://doi.org/10.1093/brain/awp301).
- 841 **Cools R**, Dopaminergic modulation of cognitive function-implications for L-DOPA treatment in Parkinson's dis-  
842 ease; 2006. doi: [10.1016/j.neubiorev.2005.03.024](https://doi.org/10.1016/j.neubiorev.2005.03.024).
- 843 **Cools R**, Gibbs SE, Miyakawa A, Jagust W, D'Esposito M. Working memory capacity predicts dopamine synthesis  
844 capacity in the human striatum. *Journal of Neuroscience*. 2008; 28. doi: [10.1523/JNEUROSCI.4475-07.2008](https://doi.org/10.1523/JNEUROSCI.4475-07.2008).
- 845 **Cools R**, Sheridan M, Jacobs E, D'Esposito M. Impulsive personality predicts dopamine-dependent changes in  
846 frontostriatal activity during component processes of working memory. *Journal of Neuroscience*. 2007; 27.  
847 doi: [10.1523/JNEUROSCI.0601-07.2007](https://doi.org/10.1523/JNEUROSCI.0601-07.2007).
- 848 **Cowan N**, Chapter 20 What are the differences between long-term, short-term, and working memory?; 2008.  
849 doi: [10.1016/S0079-6123\(07\)00020-9](https://doi.org/10.1016/S0079-6123(07)00020-9).
- 850 **Dilmore JG**, Gutkin BS, Ermentrout GB. Effects of dopaminergic modulation of persistent sodium currents on  
851 the excitability of prefrontal cortical neurons: A computational study. *Neurocomputing*. 1999; 26-27. doi:  
852 [10.1016/S0925-2312\(99\)00005-3](https://doi.org/10.1016/S0925-2312(99)00005-3).
- 853 **Durstewitz D**, Seamans JK, Sejnowski TJ. Neurocomputational Models of Working Memory. *Nature Neuro-*  
854 *science*. 2000; 3. doi: [10.1038/81460](https://doi.org/10.1038/81460).
- 855 **Edin F**, Klingberg T, Johansson P, McNab F, Tegnér J, Compte A. Mechanism for top-down control of working  
856 memory capacity. *Proceedings of the National Academy of Sciences of the United States of America*. 2009;  
857 106. doi: [10.1073/pnas.0901894106](https://doi.org/10.1073/pnas.0901894106).
- 858 **Feldmann-Wüstefeld T**, Vogel EK. Neural Evidence for the Contribution of Active Suppression During Working  
859 Memory Filtering. *Cerebral Cortex*. 2019; 29. doi: [10.1093/cercor/bhx336](https://doi.org/10.1093/cercor/bhx336).
- 860 **Fischer J**, Whitney D. Serial dependence in visual perception. *Nature Neuroscience*. 2014; 17:738–743. doi:  
861 [10.1038/nn.3689](https://doi.org/10.1038/nn.3689).
- 862 **Frank MJ**. Dynamic dopamine modulation in the basal ganglia: A neurocomputational account of cognitive  
863 deficits in medicated and nonmedicated Parkinsonism. *Journal of Cognitive Neuroscience*. 2005; 17. doi:  
864 [10.1162/0898929052880093](https://doi.org/10.1162/0898929052880093).
- 865 **Frank MJ**, Badre D, Mechanisms of hierarchical reinforcement learning in corticostriatal circuits 1: Computa-  
866 tional analysis; 2012. doi: [10.1093/cercor/bhr114](https://doi.org/10.1093/cercor/bhr114).
- 867 **Frank MJ**, Claus ED. Anatomy of a decision: Striato-orbitofrontal interactions in reinforcement learning, deci-  
868 sion making, and reversal. *Psychological Review*. 2006; 113. doi: [10.1037/0033-295X.113.2.300](https://doi.org/10.1037/0033-295X.113.2.300).
- 869 **Frank MJ**, Fossella JA, Neurogenetics and pharmacology of learning, motivation, and cognition; 2011. doi:  
870 [10.1038/npp.2010.96](https://doi.org/10.1038/npp.2010.96).
- 871 **Frank MJ**, Loughry B, O'Reilly RC, Interactions between frontal cortex and basal ganglia in working memory: A  
872 computational model; 2001. doi: [10.3758/CABN.1.2.137](https://doi.org/10.3758/CABN.1.2.137).
- 873 **Fukuda K**, Awh E, Vogel EK, Discrete capacity limits in visual working memory; 2010. doi:  
874 [10.1016/j.conb.2010.03.005](https://doi.org/10.1016/j.conb.2010.03.005).
- 875 **Gerfen CR**, Molecular effects of dopamine on striatal-projection pathways; 2000. doi: [10.1016/S1471-](https://doi.org/10.1016/S1471-1931(00)00019-7)  
876 [1931\(00\)00019-7](https://doi.org/10.1016/S1471-1931(00)00019-7).
- 877 **Gorelova NA**, Yang CR. Dopamine D1/D5 receptor activation modulates a persistent sodium current in rat  
878 prefrontal cortical neurons in vitro. *Journal of Neurophysiology*. 2000; 84. doi: [10.1152/jn.2000.84.1.75](https://doi.org/10.1152/jn.2000.84.1.75).
- 879 **Hazy TE**, Frank MJ, O'Reilly RC. Towards an executive without a homunculus: Computational models of the  
880 prefrontal cortex/basal ganglia system. *Philosophical Transactions of the Royal Society B: Biological Sciences*.  
881 2007; 362. doi: [10.1098/rstb.2007.2055](https://doi.org/10.1098/rstb.2007.2055).

- 882 **Hazy TE**, Frank MJ, O'Reilly RC. Computational Neuroscientific Models of Working Memory. . 2021; .
- 883 **Hitchcock PF**, Frank MJ. The challenge of learning adaptive mental behavior. *Journal of Psychopathology and*  
884 *Clinical Science*. 2024; .
- 885 **Ito T**, Murray JD. Multitask representations in the human cortex transform along a sensory-to-motor hierarchy.  
886 *Nature Neuroscience*. 2023; .
- 887 **Jaeger D**, Kita H, Wilson CJ. Surround inhibition among projection neurons is weak or nonexistent in the rat  
888 neostriatum. *Journal of Neurophysiology*. 1994; 72. doi: [10.1152/jn.1994.72.5.2555](https://doi.org/10.1152/jn.1994.72.5.2555).
- 889 **Jaskir A**, Frank MJ. On the normative advantages of dopamine and striatal opponency for learning and choice.  
890 *eLife*. 2023; 12. doi: [10.7554/elife.85107](https://doi.org/10.7554/elife.85107).
- 891 **Kiyonaga A**, Scimeca JM, Bliss DP, Whitney D, Serial Dependence across Perception, Attention, and Memory;  
892 2017. doi: [10.1016/j.tics.2017.04.011](https://doi.org/10.1016/j.tics.2017.04.011).
- 893 **Kriete T**, Noelle DC, Cohen JD, O'Reilly RC. Indirection and symbol-like processing in the prefrontal cortex and  
894 basal ganglia. *Proceedings of the National Academy of Sciences of the United States of America*. 2013; 110.  
895 doi: [10.1073/pnas.1303547110](https://doi.org/10.1073/pnas.1303547110).
- 896 **Krueger KA**, Dayan P. Flexible shaping: How learning in small steps helps. *Cognition*. 2009; 110. doi:  
897 [10.1016/j.cognition.2008.11.014](https://doi.org/10.1016/j.cognition.2008.11.014).
- 898 **Levitt JB**, Lewis DA, Yoshioka T, Lund JS. Topography of pyramidal neuron intrinsic connections in  
899 macaque monkey prefrontal cortex (areas 9 and 46). *Journal of Comparative Neurology*. 1993; 338. doi:  
900 [10.1002/cne.903380304](https://doi.org/10.1002/cne.903380304).
- 901 **Lisman JE**, Fellous JM, Wang XJ. A role for NMDA-receptor channels in working memory. *Nature Neuroscience*.  
902 1998; 1. doi: [10.1038/1086](https://doi.org/10.1038/1086).
- 903 **Luck SJ**, Vogel EK, Visual working memory capacity: From psychophysics and neurobiology to individual differ-  
904 ences; 2013. doi: [10.1016/j.tics.2013.06.006](https://doi.org/10.1016/j.tics.2013.06.006).
- 905 **Ma WJ**, Husain M, Bays PM, Changing concepts of working memory; 2014. doi: [10.1038/nn.3655](https://doi.org/10.1038/nn.3655).
- 906 **Maia TV**, Frank MJ, An Integrative Perspective on the Role of Dopamine in Schizophrenia; 2017. doi:  
907 [10.1016/j.biopsycho.2016.05.021](https://doi.org/10.1016/j.biopsycho.2016.05.021).
- 908 **McNab F**, Klingberg T. Prefrontal cortex and basal ganglia control access to working memory. *Nature Neuro-*  
909 *science*. 2008; 11. doi: [10.1038/nn2024](https://doi.org/10.1038/nn2024).
- 910 **Moscovitch M**, Winocur G. In: *The Frontal Cortex and Working with Memory*; 2009. doi:  
911 [10.1093/acprof:oso/9780195134971.003.0012](https://doi.org/10.1093/acprof:oso/9780195134971.003.0012).
- 912 **Moustafa AA**, Cohen MX, Sherman SJ, Frank MJ. A role for dopamine in temporal decision making and reward  
913 maximization in Parkinsonism. *Journal of Neuroscience*. 2008; 28. doi: [10.1523/JNEUROSCI.3116-08.2008](https://doi.org/10.1523/JNEUROSCI.3116-08.2008).
- 914 **Nakajima M**, Schmitt LI, Halassa MM. Prefrontal Cortex Regulates Sensory Filtering through a Basal Ganglia-  
915 to-Thalamus Pathway. *Neuron*. 2019; 103. doi: [10.1016/j.neuron.2019.05.026](https://doi.org/10.1016/j.neuron.2019.05.026).
- 916 **Nassar MR**, Helmers JC, Frank MJ. Chunking as a rational strategy for lossy data compression in visual working  
917 memory. *Psychological Review*. 2018; 125. doi: [10.1037/rev0000101](https://doi.org/10.1037/rev0000101).
- 918 **Nyberg L**, Eriksson J. Working memory: Maintenance, updating, and the realization of intentions. *Cold Spring*  
919 *Harbor Perspectives in Biology*. 2016; 8. doi: [10.1101/cshperspect.a021816](https://doi.org/10.1101/cshperspect.a021816).
- 920 **Oberauer K**. The focus of attention in working memory—from metaphors to mechanisms. *Frontiers in Human*  
921 *Neuroscience*. 2013; .
- 922 **Oberauer K**, and Simon Farrell SL, Jarrold C, Greaves M. Modeling working memory: An interference model of  
923 complex span. *Psychonomic Bulletin Review*. 2012; .
- 924 **Oberauer K**, Lin HY. An interference model of visual working memory. *Psychological Review*. 2017; .
- 925 **O'Reilly RC**, Frank MJ. Making working memory work: A computational model of learning in the prefrontal  
926 cortex and basal ganglia. *Neural Computation*. 2006; 18. doi: [10.1162/089976606775093909](https://doi.org/10.1162/089976606775093909).

- 927 **O'Reilly RC**, Munakata Y. Computational Explorations in Cognitive Neuroscience: Understanding the Mind by  
928 Simulating the Brain. The MIT Press; 2000.
- 929 **O'Reilly RC**, Munakata Y. Computational Explorations in Cognitive Neuroscience; 2019. doi: [10.7551/mit-](https://doi.org/10.7551/mitpress/2014.001.0001)  
930 [press/2014.001.0001](https://doi.org/10.7551/mitpress/2014.001.0001).
- 931 **O'Reilly RC**, Munakata Y, Frank MJ, Hazy TE, Contributors. Computational Cognitive Neuroscience. Online Book,  
932 5th Edition, URL: <https://compcogneuro.org>; 2024. <https://compcogneuro.org/book>.
- 933 **Pucak ML**, Levitt JB, Lund JS, Lewis DA. Patterns of intrinsic and associational circuitry in monkey prefrontal cor-  
934 tex. Journal of Comparative Neurology. 1996; 376. doi: [10.1002/\(SICI\)1096-9861\(19961223\)376:4<614::AID-](https://doi.org/10.1002/(SICI)1096-9861(19961223)376:4<614::AID-CNE9>3.0.CO;2-4)  
935 [CNE9>3.0.CO;2-4](https://doi.org/10.1002/(SICI)1096-9861(19961223)376:4<614::AID-CNE9>3.0.CO;2-4).
- 936 **Pusch R**, Packheiser J, Azizi AH, Sevincik CS, Rose J, Cheng S, Stüttgen MC, Güntürkün O. Working memory  
937 performance is tied to stimulus complexity. Communications Biology. 2023 12; 6. doi: [10.1038/s42003-023-](https://doi.org/10.1038/s42003-023-05486-7)  
938 [05486-7](https://doi.org/10.1038/s42003-023-05486-7).
- 939 **Rikhye RV**, Gilra A, Halassa MM. Thalamic regulation of switching between cortical representations enables  
940 cognitive flexibility. Nature Neuroscience. 2018; 21. doi: [10.1038/s41593-018-0269-z](https://doi.org/10.1038/s41593-018-0269-z).
- 941 **Shen W**, Flajolet M, Greengard P, Surmeier DJ. Dichotomous dopaminergic control of striatal synaptic plasticity.  
942 Science. 2008; 321. doi: [10.1126/science.1160575](https://doi.org/10.1126/science.1160575).
- 943 **Stocco A**, Lebiere C, Anderson JR. Conditional Routing of Information to the Cortex: A Model of the Basal  
944 Ganglia's Role in Cognitive Coordination. Psychological Review. 2010; 117. doi: [10.1037/a0019077](https://doi.org/10.1037/a0019077).
- 945 **Swan G**, Wyble B. The binding pool: A model of shared neural resources for distinct items in visual working  
946 memory. Attention, Perception, and Psychophysics. 2014; 76. doi: [10.3758/s13414-014-0633-3](https://doi.org/10.3758/s13414-014-0633-3).
- 947 **Traylor A**, Merullo J, Frank MJ, Pavlick E. Transformer Mechanisms Mimic Frontostriatal Gating Operations  
948 When Trained on Human Working Memory Tasks. . 2024; .
- 949 **Urakubo H**, Honda M, Froemke RC, Kuroda S. Requirement of an allosteric kinetics of NMDA receptors for spike  
950 timing-dependent plasticity. Journal of Neuroscience. 2008; 28. doi: [10.1523/JNEUROSCI.0303-08.2008](https://doi.org/10.1523/JNEUROSCI.0303-08.2008).
- 951 **Vogel EK**, McCollough AW, Machizawa MG. Neural measures reveal individual differences in controlling access  
952 to working memory. Nature. 2005; 438. doi: [10.1038/nature04171](https://doi.org/10.1038/nature04171).
- 953 **Wei Z**, Wang XJ, Wang DH. From distributed resources to limited slots in multiple-item working memory: A  
954 spiking network model with normalization. Journal of Neuroscience. 2012; 32. doi: [10.1523/JNEUROSCI.0735-](https://doi.org/10.1523/JNEUROSCI.0735-12.2012)  
955 [12.2012](https://doi.org/10.1523/JNEUROSCI.0735-12.2012).
- 956 **Wiecki TV**, Riedinger K, Ameln-Mayerhofer AV, Schmidt WJ, Frank MJ. A neurocomputational account of  
957 catalepsy sensitization induced by D2 receptor blockade in rats: Context dependency, extinction, and re-  
958 newal. Psychopharmacology. 2009; 204. doi: [10.1007/s00213-008-1457-4](https://doi.org/10.1007/s00213-008-1457-4).
- 959 **Wilhelm M**, Sych Y, Fomins A, Warren JLA, Lewis C, Capdevila LS, Boehringer R, Amadei EA, Grewe B, O'Connor  
960 EC, Hall BJ, Helmchen F. Striatum-projecting prefrontal cortex neurons support working memory mainte-  
961 nance. Nature Communications. 2023 12; 14. doi: [10.1038/s41467-023-42777-3](https://doi.org/10.1038/s41467-023-42777-3).
- 962 **Xiao Y**, Wang Y, Felleman DJ. A spatially organized representation of colour in macaque cortical area V2. Nature.  
963 2003; 421. doi: [10.1038/nature01372](https://doi.org/10.1038/nature01372).
- 964 **Zhang D**, Zhao H, Bai W, Tian X. Functional connectivity among multi-channel EEGs when working memory  
965 load reaches the capacity. Brain Research. 2016; 1631. doi: [10.1016/j.brainres.2015.11.036](https://doi.org/10.1016/j.brainres.2015.11.036).
- 966 **Zhang W**, Luck SJ. Discrete fixed-resolution representations in visual working memory. Nature. 2008; 453. doi:  
967 [10.1038/nature06860](https://doi.org/10.1038/nature06860).

## 968 Appendix 1

### 969 Appendix

#### 970 Neural model simulations

971 For *each frontal stripe*, the corresponding striatal gating layers consisted of 24 distributed  
972 units (12 Go and 12 NoGo) which learn the probability of obtaining a reward if the stimulus  
973 in question is gated into, or out of, its respective working memory stripe. The number of  
974 units is proportional to the number of stripes and therefore in the model iterations with  
975 4 or 8 stripes, the striatal gating layers are larger. In each module, a Gpi/Thal (globus pal-  
976 lidas internus/thalamus) unit implements a gating signal and is activated when relatively  
977 more striatal Go than NoGo units are active (subject to inhibitory competition from other  
978 GPI/Thal units that modulate gating of neighboring stripes (**O'Reilly and Frank, 2006**). Thus  
979 the GPI/Thal units summarize the contributions of multiple interacting layers that imple-  
980 ment gating among the globus pallidus, subthalamic nucleus, and thalamus as simulated in  
981 more detailed networks of a single BG circuit (**Frank and Claus, 2006**), in these larger-scale  
982 networks we abstract away from these details. For input-gating circuits, GPI/Thal activation  
983 induces maintenance of activation states in the corresponding frontal maintenance layer  
984 (**Frank et al., 2001; O'Reilly and Frank, 2006**). For output-gating circuits, the GPI/Thal acti-  
985 vation results in information flow from the frontal maintenance layer to the frontal output  
986 layer. This output layer projects to the decision circuit, such that only output-gated repre-  
987 sentations influence response selection.

988 The task in this experiment was the color visual working memory task. Each trial con-  
989 sisted of stimulus presentation, during which stimuli could be gated into corresponding PFC  
990 areas, followed by another phase in which all input stimuli were removed and the network  
991 had to rely on maintained PFC representations in order to respond. The frontal stripes for  
992 each of the stimulus dimensions could independently maintain representations of these  
993 stimulus dimensions in  $PFC_{MaintDeep}$  subject to gating signals from the BG. Initially, a "Go  
994 bias" encourages exploratory updating (and subsequent maintenance) due to novelty; these  
995 gating signals are then reinforced to the extent that the frontal representations come to be  
996 predictive of reward (**O'Reilly and Frank (2006)**). However, not all maintained  $PFC_{MaintDeep}$   
997 representations influence decision in the response circuitry, only those that are also rep-  
998 resented in  $PFC_{Out}$  due to output gating signals. Thus in a given trial, color of the current  
999 stimulus may be represented in  $PFC_{MaintDeep}$  but the output gating will dictate the ultimate  
1000 model response.

#### 1001 Neural model implementational details

The model is implemented using the Leabra framework (**O'Reilly and Munakata, 2000, 2019**),  
with the new version of the emergent neural simulation software, which contains extensive  
documentation and examples that can be run in Python or the Go language (<https://github.com/emergent>).  
All of the computational models, and the code to perform the analysis, are available and will  
be published on our github account. The PBWM network used here simulates the anatomical  
projections and physiological properties of the BG circuitry in learning, working memory  
and decision making (**Frank, 2005; O'Reilly and Frank, 2006**). Leabra uses point neurons  
with excitatory, inhibitory, and leak conductances contributing to an integrated membrane  
potential, which is then thresholded and transformed to produce a rate code output com-  
municated to other units. Dopamine in the BG modifies activity in Go and NoGo units in the  
striatum, and this modulation of activity affects both the propensity for overall gating (Go  
relative to NoGo activity), and activity-dependent plasticity that occurs during reward pre-  
diction errors (**Frank, 2005; Wiecki et al., 2009; Beeler et al., 2012; Jaskir and Frank, 2023**).

1012  
1013  
1014  
1015  
1016  
1017  
1018  
1019  
1020  
1021  
1022  
1023  
1024  
1025  
1026  
1027  
1028  
1029  
1030  
1031  
1032  
1033  
1034  
1035  
1036  
1037  
1038  
1039  
1040  
1041  
1042  
1043  
1044  
1045  
1046  
1047  
1048  
1049  
1050  
1051  
1052  
1053  
1054

Both of these functions are detailed below.

The membrane potential  $V_m$  is updated as a function of ionic conductances  $g$  with reversal (driving) potentials  $E$  according to the following differential equation:

$$C_m \frac{dV_m}{dt} = g_e(t)\bar{g}_e(E_e - V_m) + g_i(t)\bar{g}_i(E_i - V_m) + g_l(t)\bar{g}_l(E_l - V_m) + \dots \quad (8)$$

(9)

where  $C_m$  is the membrane capacitance and determines the time constant with which the voltage can change, and subscripts  $e$ ,  $l$  and  $i$  refer to excitatory, leak, and inhibitory channels respectively (and "... " refers to the possibility of adding other channels implementing neural accommodation and hysteresis). Following electrophysiological convention, the overall conductance for each channel  $c$  is decomposed into a time-varying component  $g_c(t)$  computed as a function of the dynamic state of the network, and a constant  $\bar{g}_c$  that controls the relative influence of the different conductances. The equilibrium potential can be written in a simplified form by setting the excitatory driving potential ( $E_e$ ) to 1 and the leak and inhibitory driving potentials ( $E_l$  and  $E_i$ ) of 0:

$$V_m^\infty = \frac{g_e \bar{g}_e}{g_e \bar{g}_e + g_l \bar{g}_l + g_i \bar{g}_i} \quad (10)$$

which shows that the neuron is computing a balance between excitation and the opposing forces of leak and inhibition. The excitatory net input/conductance  $g_e(t)$  is computed as the proportion of open excitatory channels as a function of sending activations times the weight values:

$$g_e(t) = \langle x_i * wi \rangle = \frac{1}{n} \sum_i x_i w_i \quad (11)$$

The inhibitory conductance is computed as described in the next section, and leak is a constant.

Activation communicated to other cells ( $y_j$ ) is a thresholded ( $\Theta$ ) sigmoidal function of the membrane potential with gain parameter  $\gamma$ :

$$y = \frac{1}{\left(1 + \frac{1}{\gamma(g_e - g_e^\ominus)_+}\right)} \quad (12)$$

where  $g_e^\ominus$  is the level of excitatory input conductance that would put the equilibrium membrane potential right at the firing threshold  $\Theta$  and depends on the level of inhibition and leak.

$$g_e^\ominus = \frac{g_i(E_i - \Theta) + g_l(E_l - \Theta)}{\Theta - E_e} \quad (13)$$

### Inhibition Within Layers

For within layer lateral inhibition, Leabra uses feedforward and feedback (FFFB) inhibition, allowing the  $g_i$  values for each neuron to be adjusted as a function of total net input to a layer (feedforward inhibition) as well as the total excitatory activity of that layer (feedback

1056  
1057  
1058  
1059  
1060  
1061  
1062  
1063  
1064  
1065  
1066  
1067  
1068  
1069  
1070  
1071  
1072  
1073  
1074  
1075  
1076  
1077  
1078  
1079  
1080  
1081  
1082  
1083  
1084  
1085  
1086  
  
1087  
1088  
1089  
1090  
1091  
1092  
1093  
1094  
1095  
1096  
1097  
1098  
1099  
1100  
1101  
1102  
1103  
1104  
1105

inhibition). The average net input (excitatory conductance) to a layer is just the average of the of each unit indexed by in the layer:

$$\langle g_e \rangle = \sum_n \frac{1}{n} g e_i$$

Similarly, the average activation is just the average of the activation values ( $y_i$ ):

$$\langle y \rangle = \sum_n \frac{1}{n} y_i$$

The overall inhibitory conductance is just the sum of the two terms (ff and fb), with an overall inhibitory gain constant factor  $G_i$ :

$$g_i(t) = G_i [\text{ff}(t) + \text{fb}(t)]$$

This  $G_i$  factor is typically the only parameter manipulated to determine overall layer activity level. The default value is 1.8 (but is reduced in the chunk layer, see below).

The feedforward (ff) term is:

$$\text{ff}(t) = \text{FF} [\langle g_e \rangle - \text{FF0}]_+$$

where FF is a constant gain factor for the feedforward component (set to 1.0 by default), and FF0 is a constant offset (set to 0.1 by default).

The feedback (fb) term is:

$$\text{fb}(t) = \text{fb}(t - 1) + dt [\text{FB}\langle y \rangle - \text{fb}(t - 1)]$$

where FB is the overall gain factor for the feedback component (0.5 default),  $dt$  is the time constant for integrating the feedback inhibition (0.7 default), and the  $t-1$  indicates the previous value of the feedback inhibition.

## Connectivity

The connectivity of the BG network is critical, and is thus summarized here (see Frank, 2006 and O'Reilly & Frank, 2006 for details and references). Unless stated otherwise, projections are fully connected (that is all units from the source region target the destination region, with a randomly initialized synaptic weight matrix). However the units in PFC, Striatum, GPiThal are all organized with columnar structure. Units in the first stripe of PFC represent one set of representations and project to a single column of Go and NoGo units in the Striatum, which in turn project to the corresponding column in GPiThal. Each Thalamic unit is reciprocally connected with the associated column in PFC. This connectivity is similar to that described by anatomical studies, in which the same cortical region that projects to the striatum is modulated by the output through the BG circuitry and Thalamus.

Dopamine units in the SNc project to the entire Striatum, but with different projections to encode the effects of D1 receptors in Go neurons and D2 receptors in NoGo neurons. With increased dopamine, Go units are excited while NoGo units are inhibited, and vice-versa with lowered dopamine levels. The particular set of units that are impacted by dopamine is determined by those receiving excitatory input from sensory cortex and PFC. Thus dopamine modulates this activity, thereby affecting the relative balance of Go vs NoGo activity in those units activated by cortex. This impact of dopamine on Go/NoGo activity levels influences both the propensity for gating (during response selection) and learning, as described next.



## Learning

Learning in the model is activity dependent and using a biologically motivated homeostatic learning rule called the eXtended Contrastive Attractor Learning (XCAL) rule. The empirical learning function (called the XCAL dWt function) approximates that observed from a highly biologically detailed computational model of the known synaptic plasticity mechanisms, by [Urakubo et al. \(2008\)](#); see [O'Reilly et al. \(2024\)](#). XCAL uses a simple piecewise-linear function, described below, that emerges from it. This XCAL dWt function resembles the BCM learning function, where weight changes are a function of presynaptic activation  $x$  and postsynaptic activation  $y$  relative to a floating threshold (approximating effects of calcium levels), and is functionally similar to contrastive Hebbian learning. The floating threshold determines the amount of activity needed to elicit LTP vs LTD.

$$\Delta w = f_{xcal}(xy, \theta_p)$$

$$f_{xcal}(xy, \theta_p) = \begin{cases} (xy - \theta_p) & \text{if } xy > \theta_p \theta_d \\ -xy(1 - \theta_d)/\theta_d, & \text{otherwise} \end{cases}$$

where  $\theta_p$  is the floating threshold and  $\theta_d = 0.1$  is a constant that determines the point where the function reverses direction (i.e., back toward zero within the weight decrease regime).

In XCAL, error-driven learning is accommodated by allowing the floating threshold to vary as a function of recent activity ([O'Reilly et al., 2024](#)). Weights are increased if activity states during the outcome are greater than their recent levels (i.e. activity states while the network is generating a response), and conversely, weights decrease if the activity levels go down relative to prior states. Thus, we can think of the recent activity levels (the threshold) as reflecting expectations which are subsequently compared to actual outcomes, with the difference (or "error") driving learning. Thus XCAL is closely related to contrastive Hebbian learning, where weight changes are determined by changes in activation. As we will see below, in PBWM and BG models, the error in gating signals is driven by changes in activation resulting from dopaminergic RPEs ([Frank, 2005](#); [O'Reilly and Frank, 2006](#); [Frank and Badre, 2012](#)).

## Striatal Learning Function

Synaptic connection weights in striatal units were trained using a reinforcement learning version of Leabra. The learning algorithm involves two phases, and is more biologically plausible than standard error backpropagation. In the *minus phase*, the network settles into activity states based on input stimuli and its synaptic weights, ultimately "choosing" a response. In the *plus phase*, the network resettles in the same manner, with the only difference being a change in simulated dopamine: an increase of SNc unit firing for positive reward prediction errors, and a decrease for negative prediction errors ([Frank, 2005](#); [O'Reilly and Frank, 2006](#)). This change in dopamine modifies Go and NoGo activity levels, and because synaptic strengths are adjusted as a function of activity levels relative to the floating threshold (previous activity levels prior to the RPE), this functionality also influences what is learned.

For the large-scale BG-PFC models used here and in [O'Reilly and Frank \(2006\)](#) some abstractions are used. Each stripe (group of units) in the Striatum layer is divided into Go vs. NoGo in an alternating fashion. The DA input from the SNc modulates these unit activations in the update phase by providing extra excitatory current to Go and extra inhibitory current to the NoGo units in proportion to the positive magnitude of the DA signal, and vice-versa

1153

1154

1155

1156

1157

1158

1159

1160

1161

1162

for negative DA magnitude. This reflects the opposing influences of DA on these neurons (*Frank, 2005; Gerfen, 2000; Shen et al., 2008*). The update phase DA signal reflects the critic system's reward prediction error (RPE) produced by gating signals (see below) – that is, if the PFC state is predictive of reward, the striatal units will be reinforced. Learning on weights into the Go/NoGo units is based on the activation delta that results from this RPE using the same XCAL dWt function defined above (which is functionally similar to contrastive Hebbian learning).

1163

### Dopamine and prediction errors in the “critic”

1164

1165

1166

1167

1168

1169

1170

1171

1172

1174

1173

1175

We used a simplified version of the critic in these simulations because they do not depend on the differences between different algorithms (e.g, temporal difference learning or “PVLV”, the algorithm used in our other BG-PFC networks (*O'Reilly and Frank, 2006; Hazy et al., 2007*)). The algorithm we used for generating dopaminergic prediction errors corresponds to a basic Rescorla-Wagner delta rule, as also reported in *Frank and Badre (2012)*. The use of a simple delta rule allowed us to confirm that simulation results do not depend on the details of the basic RL algorithm. The reward prediction error (RPE or  $\delta$ ) is the difference between the delivered reward (R) and the expected reward (V). The expected reward is calculated by a RewardPrediction unit which calculates the value of the gating actions based on the previously earned rewards.

1176

$$\delta = R - V \quad (14)$$

1177

1178

1180

1179

1181

An updated V in this RewardPrediction unit is calculated after the RPE is determined, with learning rate  $\alpha$ :

1182

$$V_{updated} = V + \alpha(\delta) \quad (15)$$

1183

1184

### GPiThal Units

1185

1186

1187

1188

1189

1190

1191

The GPiThal units provide a simplified version of the GPi/Thalamus layers abstracted from the full circuitry implemented in more basic versions of the BG circuit (*Frank and Claus, 2006, e.g.,*). They receive a net input that reflects the normalized Go - NoGo activations in the corresponding Striatum stripe:

$$\alpha_j = \left[ \frac{\sum Go - \sum NoGo}{\sum Go + \sum NoGo} \right]_+ \quad (16)$$

1192

1193

1194

1195

1196

1197

1198

1199

(where  $[\ ]_+$  indicates that only the positive part is taken; when there is more NoGo than Go, the net input is 0). This net input then drives standard Leabra point neuron activation dynamics, with FFFB inhibitory competition dynamics that cause stripes to compete for both input and output gating of PFC. This dynamic is consistent with the notion that competition/selection takes place primarily in the smaller GP/GPi areas, and not much in the much larger striatum e.g. (*Jaeger et al., 1994*). The resulting GPiThal activation then provides the gating update signal to the PFC: if the corresponding GPiThal unit is active (above a minimum threshold; .1), then active maintenance currents in the PFC are toggled.

1200

### PFC Maintenance

PFC active maintenance is supported in part by excitatory ionic conductances that are toggled by Go firing from the GPiThal layers. This is implemented with an extra excitatory ion channel in the basic  $V_m$  update equation (9). This channel has a conductance value of .5 when active. See (*Frank et al., 2001*) for further discussion of this kind of maintenance mechanism, which has been proposed by several researchers (e.g.) (*Lisman et al., 1998*;

1203

1204

1205

1206

1207

1208

1209

1210

1211

1212

1213

1214

1215

1216

1217

1218

1219

1220

1221

1222

1223

1224

1225

1226

1227

1228

1229

1230

1231

1232

1233

1234

1235

1236

1237

1238

1239

1240

1241

1242

1243

1244

1245

1246

*Dilmore et al., 1999; Gorelova and Yang, 2000; Durstewitz et al., 2000*). The first opportunity to toggle PFC maintenance occurs at the end of the first plus phase, and then again at the end of the second plus phase (third phase of settling). Thus, a complete update can be triggered by two Go's in a row, and it is almost always the case that if a Go fires the first time, it will fire the next, because Striatum firing is primarily driven by sensory inputs, which remain constant.

### Continuous Stimuli

Previous applications of PBWM consisted of discrete stimuli. For this project, the stimuli were made continuous on a ring from 0 to 360 to mimic the color wheel. The Gaussian bump width is 0.15 and there are 20 units per layer.

Stimuli are presented sequentially. During a store trial, the color and orientation are presented simultaneously and remain active for the duration of the gating decision and until the next trial begins. Each trial is 100msec long (10 Hz alpha frequency) and is organized into 4 quarters. Each quarter lasts 25 msec (40 Hz, gamma frequency). The first 3 quarters form the expectation (minus phase) and the last quarter is the outcome (plus phase). The difference between the minus and plus phase dictates learning. For further details on the cycles (updating of each neurons membrane potential) and the timing of each trial, see *O'Reilly et al.*

### Layer Sizes and Inner Mechanics

Input, chunk and output layers each have 20 neurons.

The PFC is composed of 4 layers: maintenance and output layers, each with their superficial and deep components. When sensory inputs are presented, the superficial maintenance layers always represent them transiently. Only when gated, inputs are then maintained (due to thalamocortical innervation of deep layers) over time in the absence of subsequent input. The output layers manage output gating in an analogous fashion. All information stored in WM (maintenance deep layer) is presented to the output superficial layer. Only when the corresponding basal ganglia thalamic circuit issues an output gating signal will this information propagate to the output deep layer, and thus be accessible for "read out" by the hidden and output layers of the network.

Each PFC layer is an array of neurons that has 20 units for each stripe. For most of the experiments here, we used 2 stripes. Each stripe has 20 neurons to match with input/output layers. (This model could expanded where the representations in PFC would be distributed, as in some other applications of PBWM).

The control input layer (which represents the orientation and the store/recall instruction) has size based on the number of orientations. For 4 discrete orientations, the size of this layer is  $4 + 1$  (for ignore stimulus) +  $1$  (for recall) = 6 units. The ignore stimulus is not discussed in this paper, but allow for distractors to be presented, where the model would have to learn to ignore (not update) those stimuli. Early qualitative results suggest that findings hold even when the model is presented with stimuli that it must ignore. During a recall trial, the model input is activity in this recall unit in the control input layer, along with the orientation that should be recalled. No color is presented at this time.

The Striatum is compose of Matrix Go (D1-containing) and Matrix NoGo (D2-containing) units. Each matrix layer is further broken down into 2 sub-layers which represent input and output gating decisions. Each of these sublayers is an array of units of size number of stripes x size of control input layer. For 4 orientations and 2 stripes, each sublayer has  $2 \times 6$  (see previous paragraph) = 12 units. This means for both input and output gating, each of the Matrix Go and NoGo layers have 24 units. The matrix learns to perform a gating operation

1250

1251

1252

1253

1254

1255

1256

1257

1258

1259

1260

1261

1262

1263

1264

1265

1266

1267

1268

1269

1270

1271

1272

1273

1274

1275

1276

1277

1278

1279

1280

1281

1282

1283

1284

1285

1286

1287

1288

1289

1290

1291

1292

1293

1294

1295

1296

to the corresponding PFC stripe based on the given control input (e.g., store orientation 1; recall orientation 3, etc).

The globus pallidus externus (GPe) layer has one unit for each input and output stripe; for 2 stripes, the size of this layer is 4. These units receive inputs from corresponding Matrix NoGo layers, preventing gating of the corresponding input and output stripes when the NoGo units are active (i.e., when they have learned that these gating operations are not useful).

The globus pallidus internus (Gpi)/ Thalamus ( dorsomedial thalamus) layer receives projections from the Matrix Go layer to induce gating, but competing inputs from the GPe to prevent gating. The relative balance of Go vs NoGo signals for a given stripe thus determine, along with inhibitory competition within this layer, whether a given PFC stripe is gated. This implementation abstracts over the details of thalamic representations, combining GPI and Thalamus into one functional layer *O'Reilly and Frank (2006)* but see more detailed implementations of this functionality *Frank and Claus (2006); Collins and Frank (2013)*.

The reinforcement learning component is broken into 3 single unit layers, as described above: the Reward (which indicates the reward the model receives), the RewPred (which is the model's predicted reward), and the SNc (substantia nigra pars compacta) dopaminergic cells which calculates the difference between the two and generates the reward prediction error (RPE).

## Hyperparameter Search

The PBWM is a well established framework with many parameters. To maintain consistency with prior work and avoid an overly complex parameter search, default parameters were used with the following exceptions. A hyperparameter search was optimized for baseline model performance for parameters most relevant to this project: chunking layer inhibition (1.4), chunking layer gain (5), relative weight scales (connectivity strength) between the chunk layer and PFC layer (0.8), input layer to the chunk layer (0.7), and PFC layer to chunk layer (0.2; this lower value ensures that the chunk layer primarily reflects the input and is only influence by PFC if the nearest neighbor overlaps with the input). The striatal learning rate was set to default of 0.04 for set size 2 and 4, but changed to 0.06 for set size 3, which improved performance for the OCV analysis in that setsize (but overall performance is not substantially altered)). The relative impact of striatal NoGo vs Go activity on GpiThal layer for inducing a gating signal was set to 1.25. Finally to accommodate continuous representations with large Gaussian bump widths in PFC we increased the PFC gain parameter to 5 (otherwise the PFC would not maintain units with lesser activation on the tails of the continuous distribution). The reward function was also altered due to the continuous nature of the outputs, such that rewards are determined based on decoded values of the output layer across the population code relative to the true color that should have been reported, as a continuous linear function of the error. We also varied dopamine burst and dip gain across a wide range to explore its impact as described in the main text.

## Neural Model Output Analysis

### Out of Cluster Variance (OCV) Analysis

The OCV analysis was performed using set size of 3, largely for interpretability. (This analysis becomes convoluted when increasing the set size: the combinations of possible chunking make it hard to divide the stimuli into "other stimuli" and "probed stimuli".)

The OCV analysis was performed on a subset of trials to specifically test whether the network can leverage chunking abilities to free up resources for other items outside of the chunk. In these trials the "chunkable" items were presented first and were within 20 degrees

1297

1298

1299

1300

1301

1302

1303

1304

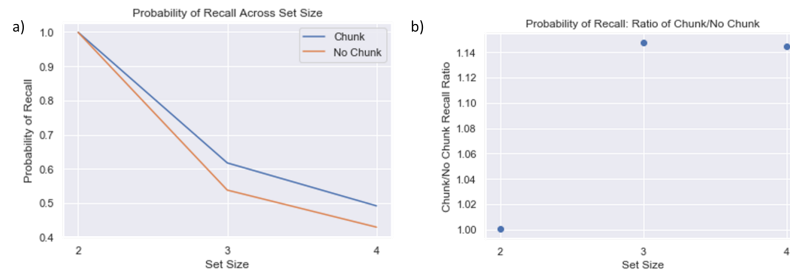
1305

1306

1307

away from each other (in color space). The third item is at least 50 degrees away from one of the stimuli. The same procedure was applied to the no-chunk control. If the network chunks the first two items it should then be more likely to store and recall the third item.

The OCV of a single trial is the variance (across the color dimension) of all stimuli in the trial except for the probed item (*Nassar et al., 2018*). The OCV is calculated for each trial and plotted against error on that trial. Since the stimuli and trials are randomly generated, some parts of the graph will be more densely populated. To combat this issue, the trials were binned so that each bin has an equal number of trials and the graph has 20 bins.



1308

**Figure 5—figure supplement 1. P(Recall) Across Set Size** a) Average recall probability across set sizes decreases with set size, but less so for chunk models. Note that chance performance is approximately 19%. b) Chunk models have a higher ratio of recall probability relative to no chunk model when set size exceeds allocated capacity. This analysis includes trials where the variance across colors is low (standard deviation was < 35 degrees). The same chunk advantages would occur across all trials (including with high variance; not shown) but we focus on low variance trials wherein models can perform reasonably well even if accidentally mistaking one item for another (swap errors). Here we confirm that the chunk model improvement occurs over and above such effects.