

# 1 Working memory load strengthens 2 reward prediction errors.

---

3 Abbreviated title: Working memory load strengthens prediction errors

4

5 Authors: Anne G.E. Collins<sup>1-3</sup>, Brittany Ciullo<sup>3</sup>, Michael J Frank<sup>3-4</sup>, David Badre<sup>3-4</sup>

6

7 Affiliations

8 1. Department of Psychology, University of California, Berkeley, Berkeley, CA  
9 94720.

10 2. Helen Wills Neuroscience Institute, Berkeley, CA, 94720

11 3. Department of Cognitive, Linguistics and Psychological Sciences, Brown  
12 University, Providence RI 02912

13 4. Brown Institute for Brain Science, Providence, RI 02912.

14

15 Corresponding author:

16 Anne G.E. Collins

17 3210 Tolman Hall, Department of Psychology, UC Berkeley, Berkeley, CA 94720

18 [annecollins@berkeley.edu](mailto:annecollins@berkeley.edu),

19

20 Number of pages: 30

21 Number of figures: 6

22 Number of tables: 3

23 Abstract: 185 words

24 Introduction: 497 words

25 Discussion: 1487

26

27

28 The authors declare no competing financial interests.

29

30 Acknowledgements: We thank Christopher R Gagne for his role in data collection.

31

## 32 **Abstract:**

33 Reinforcement learning in simple instrumental tasks is usually modeled as a monolithic  
34 process in which reward prediction errors are used to update expected values of choice  
35 options. This modeling ignores the different contributions of different memory and  
36 decision-making systems thought to contribute even to simple learning. In an fMRI  
37 experiment, we asked how working memory and incremental reinforcement learning  
38 processes interact to guide human learning. Working memory load was manipulated by  
39 varying the number of stimuli to be learned across blocks. Behavioral results and  
40 computational modeling confirmed that learning was best explained as a mixture of two  
41 mechanisms: a fast, capacity-limited, and delay-sensitive working memory process  
42 together with slower reinforcement learning. Model-based analysis of fMRI data showed  
43 that striatum and lateral prefrontal cortex were sensitive to reward prediction error, as  
44 shown previously, but critically, these signals were reduced when the learning problem  
45 was within capacity of working memory. The degree of this neural interaction related to  
46 individual differences in the use of working memory to guide behavioral learning. These  
47 results indicate that the two systems do not process information independently, but  
48 rather interact during learning.

## 49 **Significance Statement**

50 Reinforcement learning theory has been remarkably productive at improving our  
51 understanding of instrumental learning as well as dopaminergic and striatal network  
52 function across many mammalian species. However, this neural network is only one  
53 contributor to human learning, and other mechanisms such as prefrontal cortex working  
54 memory, also play a key role. Our results show in addition that these other players  
55 interact with the dopaminergic RL system, interfering with its key computation of reward  
56 predictions errors.

## 57 **Intro:**

58 Reinforcement learning (RL) theory (Sutton & Barto 1998) proposes that we can learn  
59 the value associated with various choices by computing the discrepancy between the  
60 reward we obtain and our previously estimated value, and proportionally adjusting our  
61 estimate. This discrepancy, the reward prediction error (RPE), signals a valenced

62 surprise at the outcome being better or worse than expected and a direction to adapt  
63 behavior (Pessiglione et al. 2006; Schönberg et al. 2007; Daw & Doya 2006). In the  
64 brain, cortico-basal ganglia loops appear to implement a form of algorithmic RL:  
65 Dopamine-dependent plasticity in the striatum may reinforce selection of choices leading  
66 to positive RPEs and weaken those leading to negative RPEs (Frank et al. 2004; Collins  
67 & Frank n.d.). Dopaminergic neurons exhibit phasic changes in their spike rates that  
68 convey RPEs (Montague et al. 1996; Schultz 2002), and dopamine release in target  
69 regions provides a bidirectional RPE signal (Hart et al. 2014). Human imaging studies  
70 have indeed found that striatal BOLD correlates with RPE and is enhanced by DA  
71 manipulations (Pessiglione et al. 2006; Schönberg et al. 2007; Jocham et al. 2011).

72

73 However, other neurocognitive processes contribute to learning besides the integration  
74 of reward history by RL. Specifically, executive processes (such as those involved in  
75 representing sequential or hierarchical task structure) contribute substantially to human  
76 learning over and above incremental RL (Daw et al. 2011; Badre & Frank 2011;  
77 Botvinick et al. 2009; Collins & Koechlin 2012; Collins & Frank 2013). Even in basic  
78 stimulus-response learning tasks, working memory (WM) contributes substantially to  
79 instrumental learning beyond RL (Collins & Frank 2012; Collins et al. 2014), as  
80 evidenced by both behavioral analyses and quantitative computational model fits. Two  
81 effects of WM were evident in learning. As WM set size increased (working memory  
82 load), learning curves per stimulus were slowed. Second, accuracy per trial declined as  
83 a function of the number of intervening items (working memory delay). These WM  
84 effects decayed with further experience, as the more reliable but slower RL process  
85 gained control of behavior. A hybrid model of WM and RL provided a better fit to these  
86 data than either process itself (Collins & Frank 2012; Collins et al. 2014),.

87

88 This prior behavioral work implies that WM contributes to RL processes. Here, we  
89 investigate the neural markers of learning and RPEs to determine whether they are  
90 interact with WM. While many RL studies have revealed neural correlates of RPEs that  
91 relate to learning, these studies have not manipulated or estimated WM factors that  
92 could contribute to (and potentially confound) these signals. Identifying separate markers  
93 of systems that contribute jointly to behavior also provides an opportunity to explore  
94 whether they interact (e.g., competitively or cooperatively). Specifically, we tested  
95 whether frontoparietal networks associated with cognitive control and striatal systems

96 associated with RL would show parametric modulations of RPE signaling as a function  
97 of WM load during learning. We also tested whether such interactions would be  
98 predictive of the extent to which individuals relied on WM contributions to RL  
99 behaviorally.

## 100 **Methods:**

### 101 **Participants:**

102 We scanned 26 participants (ages 18-31, mean age 23, 15 males/11 females). All 26  
103 participants are included in the behavioral analysis. 5 participants were excluded from  
104 fMRI analysis prior to analyzing their fMRI data due to head movement greater than our  
105 voxel size. 2-6 blocks were excluded from 3 other participants due to movement during  
106 data collection towards the end of the scan. All participants were right-handed with  
107 normal or corrected-to-normal vision and were screened for the presence of psychiatric  
108 or neurological conditions and contraindications for fMRI. All participants were  
109 compensated for their participation and gave informed, written consent as approved by  
110 the Human Research Protection Office of Brown University.

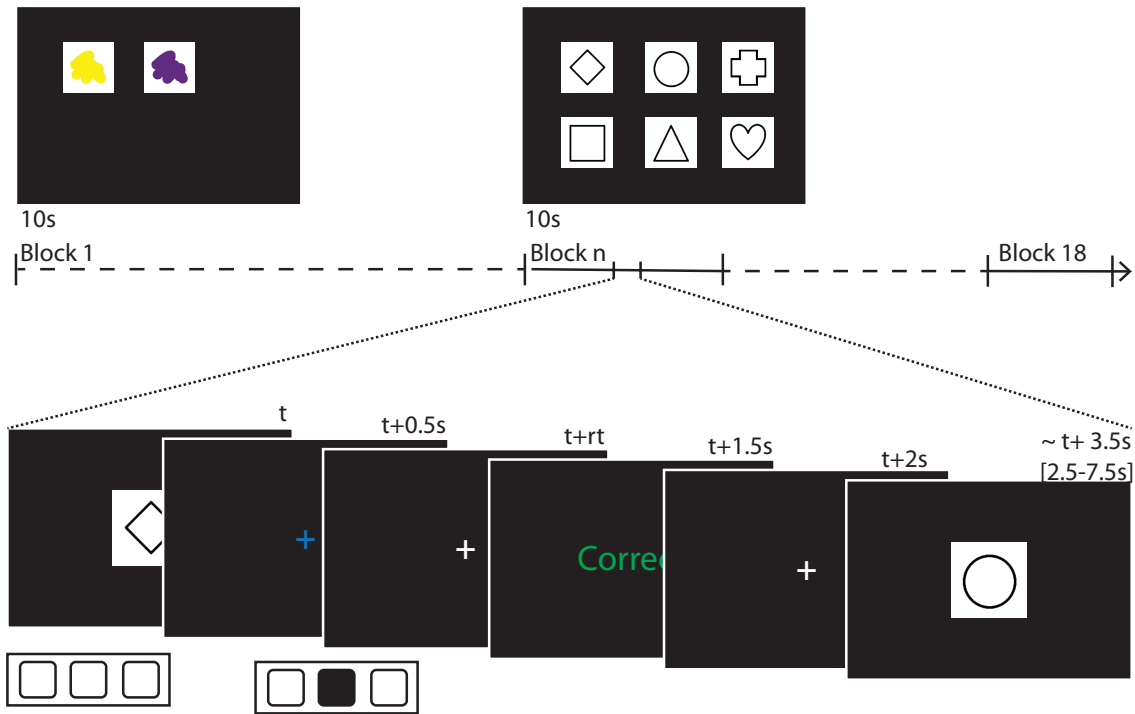
### 111 **Experimental design:**

112 The task (Fig1) was similar to that described previously (Collins & Frank 2012; Collins et  
113 al. 2014),, itself adapted from a classic Conditional Associative Learning paradigm  
114 (Petrides 1985). On each trial, subjects had to respond with one of three responses  
115 (button presses on a response pad) when presented with a centrally displayed single  
116 stimulus. Subjects had to learn over trials which response was correct for each stimulus,  
117 based on binary deterministic reinforcement feedback (Collins & Frank 2012; Collins et  
118 al. 2014),.

119

120 To manipulate working memory demands separately from RL components, we  
121 systematically varied the number of stimuli (denoted as set size  $n_s$ ) to be learned within  
122 a block. Larger set sizes provide greater load on working memory, and also impose on  
123 average larger delays between repetitions of the same stimulus. Subjects experienced 3  
124 blocks of each of the set-sizes one through six. In each block, subjects learned about a  
125 different category of visual stimulus (such as sports, fruits, places, etc.), with stimulus

126 category assignment to block set size counterbalanced across subjects. Block ordering  
 127 was also counterbalanced within subjects to ensure an even distribution of high/low load  
 128 blocks across each third of the experiment.



129  
 130 **Figure 1: Experimental Protocol.** At the beginning of each block, subjects were shown for 10 s  
 131 the set of stimuli they would see in that block. In this example, Block 1 uses color patches for  
 132 stimuli and has a set size  $n_s = 2$ ; Block n uses shapes and has  $n_s = 6$ . Each trial included the  
 133 presentation of a stimulus for 0.5s followed by a blue fixation cross until subject pressed one of  
 134 three buttons, or up to 1.5s after trial onset. Button press caused the fixation cross to turn white.  
 135 Feedback was presented for 1s, and came 1.5s after trial onset. Feedback consisted of the words  
 136 correct or incorrect in green and red, respectively. The inter-trial interval consisted of a white  
 137 fixation cross with jittered duration to allow trial by trial event-related analysis of fMRI signal.  
 138 Blocks set sizes varied between 1 and 6, and the order was randomized across subjects.

139  
 140 At the beginning of each block, subjects were shown the entire set of stimuli for that  
 141 block and were encouraged to familiarize themselves with them for a duration of 10 sec  
 142 (figure 1 top). They were then asked to make their response as quickly and accurately  
 143 as possible after each individual stimulus presentation. Within each block, stimuli were  
 144 presented 12 times each in a pseudo-randomly intermixed order.

145  
 146 Stimuli were presented in the center of the screen for up to 0.5s seconds, followed by a  
 147 blue fixation cross for up to 1s or subjects making a choice by pressing one of 3 buttons,  
 148 at which time the fixation cross turned white (figure 1 bottom). Feedback was presented  
 149 1.5s after stimulus onset for 0.5s as either “Correct” in green, “Incorrect” in red, or “Too

150 slow” if the subject failed to answer within 1.5s. A white fixation cross followed with  
151 jittered duration of mean 1.5s, ranging [1.5 6.5]s, before the next stimulus was presented.

152

153 Subjects were instructed that finding the correct action for one stimulus was not  
154 informative about the correct action for another stimulus. This was enforced in the choice  
155 of correct actions, such that, in a block with e.g.,  $n_S=3$ , the correct actions for the three  
156 stimuli were not necessarily three distinct keys. This procedure was implemented to  
157 ensure independent learning of all stimuli (i.e., to prevent subjects from inferring the  
158 correct actions to stimuli based on knowing the actions for other stimuli). Prior to  
159 entering the scanner, subjects went through the instructions and practiced on a separate  
160 set-size 2 sets of images to ensure they were familiarized with the task.

### 161 **Computational model:**

#### 162 **RLWM model:**

163 To better account for subjects’ behavior and disentangle roles of working memory and  
164 reinforcement learning, we fitted subjects’ choices with our hybrid RLWM computational  
165 model. Previous research showed that this model, allowing choice to be a mixture  
166 between a classic delta rule reinforcement learning process and a fast but capacity-  
167 limited and delay-sensitive working memory process, provided a better quantitative fit to  
168 learning data than models of either WM or RL alone (Collins & Frank 2012; Collins et al.  
169 2014). The model used here is a variant of the previously published models. We first  
170 summarize its key properties, following by the details:

- 171 • RLWM includes two modules which separately learn the value of stimulus-response  
172 mappings: a standard incremental procedural RL module with learning rate  $\alpha$ , and a  
173 WM module that updates S-R-O associations in a single trial (learning rate 1) but is  
174 capacity-limited (with capacity  $K$ ).
- 175 • The final action choice is determined as a weighted average over the two modules’  
176 policies. How much weight is given to WM relative to RL (the mixture parameter) is  
177 dynamic and reflects the probability that a subject would use WM vs. RL in guiding  
178 their choice. This weight depends on two factors. First, a *constraint* factor reflects the  
179 *a priori* probability that the item is stored in WM, which depends on set size  $n_S$  of the  
180 current block relative to capacity  $K$  (i.e., if  $n_S > K$ , the probability that an item is stored  
181 is  $K/n_S$ ), scaled by the subject’s overall reliance of WM vs. RL (factor  $0 < \rho < 1$ ), with  
182 higher values reflecting relative greater confidence in WM function. Thus, the

183 constraint factors indicates that the maximal use of WM policy relative to RL policy is  
 184  $w_0 = \rho \times \min(1, K/n_S)$ . Second, a *strategic* factor reflects the inferred reliability of the  
 185 WM compared to RL modules over time: initially, the WM module is more successful  
 186 at predicting outcomes than the RL module, but because it has higher capacity and  
 187 less vulnerability to delay, the RL module becomes more reliable with experience.

- 188 • Both RL and WM modules are subject to forgetting (decay parameters  $\phi_{RL}$  and  $\phi_{WM}$ ).

189 We constrain  $\phi_{RL} < \phi_{WM}$  consistent with WM's dependence on active memory).

190  
 191

192 **Learning model details.**

193 **Reinforcement learning model:** All models include a standard RL module with simple  
 194 delta rule learning. For each stimulus  $s$ , and action  $a$ , the expected reward  $Q(s,a)$  is  
 195 learned as a function of reinforcement history. Specifically, the  $Q$  value for the selected  
 196 action given the stimulus is updated upon observing each trial's reward outcome  $r_t$  (1 for  
 197 correct, 0 for incorrect) as a function of the prediction error between expected and  
 198 observed reward at trial  $t$ :

199 
$$Q_{t+1}(s,a) = Q_t(s,a) + \alpha \times \delta_t,$$

200 where  $\delta_t = r_t - Q_t(s,a)$  is the prediction error, and  $\alpha$  is the learning rate. Choices are  
 201 generated probabilistically with greater likelihood of selecting actions that have higher  $Q$   
 202 values, using the softmax choice rule:

203 
$$p(a|s) = \exp(\beta Q(s,a)) / \sum_i (\exp(\beta Q(s,a_i))).$$

204 Here,  $\beta$  is an inverse temperature determining the degree with which differences in  $Q$ -  
 205 values are translated into more deterministic choice, and the sum is over the three  
 206 possible actions  $a_i$ .

207

208 **Undirected noise.** The softmax temperature allows for stochasticity in choice, but where  
 209 stochasticity is more impactful when the value of actions are similar to each other. We  
 210 also allow for “slips” of action (“irreducible noise”, i.e., even when  $Q$  value differences  
 211 are large). Given a model's policy  $\pi = p(a|s)$ , adding undirected noise consists in  
 212 defining the new mixture policy:

213 
$$\pi' = (1 - \epsilon) \pi + \epsilon U,$$

214 where  $U$  is the uniform random policy ( $U(a) = 1/n_A$ ,  $n_A=3$ ), and the parameter  $0 < \epsilon < 1$   
 215 controls the amount of noise (Collins & Koechlin 2012; Collins & Frank 2013; Guitart-

216 Masip et al. 2012). (Nassar & Frank 2016) showed that failing to take into account this  
217 irreducible noise can render fits to be unduly influenced by rare odd datapoints (e.g. that  
218 might arise from attentional lapses), and that this problem is remedied by using a hybrid  
219 softmax- $\epsilon$ -greedy choice function as used here.

220

221 **Forgetting.** We allow for potential decay or forgetting in Q-values on each trial,  
222 additionally updating all Q-values at each trial, according to:

$$223 \quad Q \leftarrow Q + \phi (Q_0 - Q),$$

224 where  $0 < \phi < 1$  is a decay parameter pulling at each trial the estimates of values towards  
225 initial value  $Q_0 = 1/n_A$ . This parameter allows us to capture delay-sensitive aspects of  
226 WM, where active maintenance is increasingly likely to fail with intervening time and  
227 other stimuli, but also allows us to separately estimate any decay in RL values (which is  
228 typically substantially lower than in WM).

229

230 **Perseveration.** To allow for potential neglect of negative, as opposed to positive  
231 feedback, we estimate a perseveration parameter *pers* such that for negative prediction  
232 errors ( $\delta < 0$ ), the learning rate  $\alpha$  is reduced by  $\alpha = (1 - \textit{pers}) \times \alpha$ . Thus, values of *pers*  
233 near 1 indicate perseveration with complete neglect of negative feedback, whereas  
234 values near 0 indicate equal learning from negative and positive feedback.

235

236 **Working Memory.** To implement an approximation of a rapid updating but capacity-  
237 limited WM, this module assumes a learning rate  $\alpha = 1$  (representing the immediate  
238 accessibility of items in active memory), but includes capacity limitation such that only at  
239 most  $K$  stimuli can be remembered. At any trial, the probability of working memory  
240 contributing to the choice for a given stimulus is  $w_{WM}(t) = P_i(WM)$ . This value is dynamic  
241 as a function of experience (see next paragraph). As such, the overall policy is:

$$242 \quad \pi = w_{WM}(t)\pi_{WM} + (1 - w_{WM}(t))\pi_{RL}$$

243 where  $\pi_{WM}$  is the WM softmax policy, and  $\pi_{RL}$  is the RL policy. Note that this  
244 implementation assumes that information stored for each stimulus in working memory  
245 pertains to action-outcome associations. Furthermore, this implementation is an  
246 approximation of a capacity/resource-limited notion of working memory. It captures key  
247 aspects of working memory such as 1) rapid and accurate encoding of information when  
248 low amount of information is to be stored; 2) decrease in the likelihood of storing or  
249 maintaining items when more information is presented or when distractors are presented



250 during the maintenance period; 3) decay due to forgetting. Because it is a probabilistic  
 251 model of WM, it cannot capture specifically which items are stored, but it can provide the  
 252 likelihood of any item being accessible during choice given the task structure and recent  
 253 history (set size, delay, etc).

254

255 **Inference:** The weighting of whether to rely more on WM vs. RL is dynamically adjusted  
 256 over trials within a block based on which module is more likely to predict correct  
 257 outcomes. The initial probability of using WM  $w_{WM}(0) = P_0(WM)$  is initialized by the a  
 258 priori use of WM, as defined above,  $w_{WM}(0) = \rho \times \min(1, K/n_S)$ , where  $\rho$  is a free  
 259 parameter representing the participant's overall reliance on WM over RL.

260 On each correct trial,  $w_{WM}(t) = P_t(WM)$  is updated based on the relative likelihood that  
 261 each module would have predicted the observed outcome given the selected correct  
 262 action  $a_c$ ; specifically:

- 263 - for WM,  $p(\text{correct}|\text{stim}, WM) = w_{WM} \pi_{WM}(a_c) + (1-w_{WM})1/n_A$
- 264 - for RL,  $p(\text{correct}|\text{stim}, RL)$  this is simply  $\pi_{RL}(a_c)$

265 The mixture weight is updated by computing the posterior using the previous trial's prior,  
 266 and the above likelihoods, such that

$$P_{t+1}(WM) = \frac{P_t(WM) \times p(\text{correct}|\text{stim}, WM)}{P_t(WM) \times p(\text{correct}|\text{stim}, WM) + P_t(RL) \times p(\text{correct}|\text{stim}, RL)}$$

267 and  $P_{t+1}(RL) = 1 - P_{t+1}(WM)$ .

268

269

270 **Models Considered.** We combined the previously described features into different  
 271 learning models and conducted extensive comparisons of multiple models to determine  
 272 which fit the data best (penalizing for complexity) so as to validate the use of this model  
 273 in interpreting subjects' data. For all models we considered, adding undirected noise,  
 274 forgetting and perseveration features significantly improved the fit, accounting for added  
 275 model complexity (see model comparisons).

276

277 This left three relevant classes of models to consider:

- 278 - RL: This model combines the basic delta rule RL with forgetting, perseveration  
 279 and undirected noise features. It assumes a single system that is sensitive to  
 280 delay and asymmetry in feedback processing. This is a 5-parameter model

- 281 (learning rate  $\alpha$ , softmax inverse temperature  $\beta$ , undirected noise  $\epsilon$ , decay  $\phi_{RL}$ ,  
282 and *pers* parameter).
- 283 - RL6: This model is identical to the previous one, with the variant that learning  
284 rate can vary as a function of set-size. We have previously shown that while such  
285 a model can capture the basic differences in learning curves across set-sizes by  
286 fitting lower learning rates with higher set sizes, it provides no mechanism that  
287 would explain these effects, and still cannot capture other more nuanced effects  
288 (e.g. changes in the sensitivity to delay with experience). However it provides a  
289 benchmark to compare with RLWM. This is a 10-parameter model (6 learning  
290 rate  $\alpha_{ns}$ , softmax inverse temperature  $\beta$ , undirected noise  $\epsilon$ , decay  $\phi_{RL}$ , and *pers*  
291 parameter).
- 292 - RLWM: This is the main model, consisting of a hybrid between RL and WM. RL  
293 and WM modules have shared softmax  $\beta$  and *pers* parameters, but separate  
294 decay parameters,  $\phi_{RL}$  and  $\phi_{WM}$ , to capture their differential sensitivity to delay.  
295 Working memory capacity is  $0 < K < 6$ , with an additional parameter for overall  
296 reliance on working memory  $0 < \rho < 1$ . Undirected noise is added to the RLWM  
297 mixture policy. This is an 8-parameter model (capacity  $K$ , WM reliance  $\rho$ , WM  
298 decay  $\phi_{WM}$ , RL learning rate  $\alpha$ , RL decay  $\phi_{RL}$ , softmax inverse temperature  $\beta$ ,  
299 undirected noise  $\epsilon$ , and *pers* parameter).

300

301 In the RLWM model presented here, the RL and WM modules are independent, and only  
302 compete for choice at the policy level. Given our findings showing an interaction  
303 between the two processes, we also considered variants of RLWM including  
304 mechanisms for interactions between the two processes at the learning stage. These  
305 models provided similar fit (measured by AIC) to the simpler RLWM model. We chose to  
306 use the simpler RLWM model, because the more complex model is less identifiable  
307 within this experimental design, providing less reliable parameter estimates and  
308 regressors for model-based analysis.

309

310 **RLWM fitting procedure:** We used matlab optimization under constraint function  
311 `fmincon` to fit parameters. This was iterated with 50 randomly chosen starting points, to  
312 increase likelihood of finding a global rather than local optimum. For models including

313 the discrete capacity  $K$  parameter, this fitting was performed iteratively for capacities  $K =$   
314  $\{1,2,3,4,5\}$ , using the value gave the best fit in combination with other parameters.

315

316 Softmax  $\beta$  temperature was fit with constraints  $[0\ 100]$ . All other parameters were fit with  
317 constraints  $[0\ 1]$ . We considered sigmoid-transforming the parameters to avoid  
318 constraints in optimization and obtain normal distributions, but while fit results were  
319 similar, distributions obtained were actually not normal. Thus, all statistical tests on  
320 parameters were non-parametric. See table 4 for fit parameter statistics.

321

### 322 **Other competing models:**

323 In order to further test whether “single system” models, as opposed to hybrid models  
324 including an RL and a WM component, could account for behavior, we tested other  
325 algorithms embodying alternative assumptions in which behavior is governed by a single  
326 learning process (either RL or WM).

327 - The WMd model is similar to a WM module, with the following changes. A) there  
328 is no capacity limitation. B) Instead of being fixed, the decay parameter is fixed to  
329 an initial value which then decreases toward 0 with each stimulus encounter,  
330 modeling the possibility that forgetting in WM itself might decrease with practice.  
331 This model includes 5 parameters:  $\beta$ ,  $\varepsilon$  and  $pers$  as defined above, the initial  
332 value of decay  $decay_0$ , and  $\xi$  the decay factor.

333 - The WMdi model adds an interference mechanism to WMd, such that the decay  
334 factor of a given stimulus additionally increases with every encounter of a  
335 different stimulus. This adds one parameter to the previous model.

336 - The RL<sub>i</sub> model is identical to the basic RL model, with an added interference  
337 mechanism: on each trial, the Q-value of non-observed stimuli with the chosen  
338 action is updated in the same way as the observed stimuli, but with a fraction of  
339 the learning rate  $\alpha_i$ . This captures the possibility that credit is assigned to the  
340 wrong stimulus, modeling the possibility that WM-like effects might reflect  
341 interference within a pure RL system. This model includes 6 parameters.

342

343 **Model Comparison:** We used the Akaike Information Criterion to penalize model  
344 complexity - AIC (Burnham & Anderson 2002). Indeed, we previously showed that in the  
345 case of the RLWM model and its variants, AIC was a better approximation than  
346 Bayesian Information Criterion (BIC; Schwarz 1978) at recovering the true model from

347 generative simulations (Collins & Frank 2012). Comparing RLWM, RL6 and RL-only  
348 showed that models RL6 and RL-only were strongly non-favored, with probability 0 over  
349 the whole group. Other single process models were also unable to capture behavior  
350 better than RLWM (Fig. 3 E).

351

352 **Model Simulation:** Model selection alone is insufficient to assess whether the best  
353 fitting model sufficiently captures the data. To test whether models capture the key  
354 features of the behavior (e.g., learning curves), we simulated each model with fit  
355 parameters for each subject, with 100 repetitions per subject then averaged to represent  
356 this subject's contribution. In order to account for initial biases, we assume that the  
357 model's choice at first encounter of a stimulus is identical to the subjects, while all further  
358 choices are randomly selected from the model's learned values and policies.

359

#### 360 **fMRI recording and preprocessing:**

361 Whole-brain imaging was performed on a Siemens 3T TIM Trio MRI system equipped  
362 with a 32-channel head coil. A high-resolution T1-weighted 3D multi-echo MPRAGE  
363 image was collected from each participant for anatomical visualization. Functional  
364 images were acquired in one run of 1,920 volume acquisitions using a gradient-echo,  
365 echo planar pulse sequence (TR 2 s, TE 28 ms, flip angle 90, 40 interleaved axial slices,  
366 192 mm field of view with 3x3x3 mm voxel size). Stimuli were presented on a BOLD  
367 screen display device ([http://www.crsf.com/tools-for-functional-imaging/mr-safe-](http://www.crsf.com/tools-for-functional-imaging/mr-safe-displays/boldscreen-24-lcd-for-fmri/)  
368 [displays/boldscreen-24-lcd-for-fmri/](http://www.crsf.com/tools-for-functional-imaging/mr-safe-displays/boldscreen-24-lcd-for-fmri/)) located behind the scanner and made visible to the  
369 participant via an angled mirror attached to the head coil. Padding around the head was  
370 used to restrict head motion. Participants made their responses using an MRI-  
371 compatible button box.

372 Functional images were preprocessed in SPM8 (<http://www.fil.ion.ucl.ac.uk/spm>). Before  
373 preprocessing, data were inspected for artifacts and excessive variance in global signal  
374 (functions: `tsdiffana`, `art_global`, `art_movie`). Functional data were corrected for  
375 differences in slice acquisition timing by resampling slices to match the first slice. Next,  
376 functional data were realigned (corrected for motion) using B-spline interpolation and  
377 referenced to the mean functional image. Functional and structural images were  
378 normalized to Montreal Neurological Institute (MNI) stereotaxic space using affine  
379 regularization followed by a nonlinear transformation based on a cosine basis set, and

380 then resampled into 2x2x2 mm voxels using trilinear interpolation. Lastly, images were  
381 spatially smoothed with an 8 mm full-width at half-maximum isotropic Gaussian kernel.

## 382 **GLMs:**

383 A temporal high-pass filter of 400 seconds (.0025 Hz) was applied to our functional data  
384 in order to remove noise but preserve power from low-frequency regressors. Changes in  
385 MR signal were modeled using a general linear model (GLM) approach. Our GLM  
386 included six onsets regressors, one for correct trials corresponding to each set size (1-  
387 6). Each onset was coded as a boxcar of 2 seconds in length that encompasses  
388 stimulus presentation, response, and feedback. Each onset regressor was modulated by  
389 a Prediction Error parametric regressor. We modeled Error trials, No Response trials,  
390 and instructions (1 instruction screen at the beginning of each block, 18 total and each  
391 10 seconds in length) as separate regressors. Note that error trials across all set sizes  
392 were binned into one regressor due to the low number of error trials in low set sizes.  
393 Finally, we included nuisance regressors for the six motion parameters (x, y, z, pitch,  
394 roll, yaw) and a linear drift over the course of the run. SPM-generated regressors were  
395 created by convolving onset boxcars and parametric functions with the canonical  
396 hemodynamic response (HRF) function and the temporal derivative of the HRF. Beta  
397 weights for each regressor were estimated in a first-level, subject-specific fixed-effects  
398 model. For group analysis, the subject-specific beta estimates were analyzed with  
399 subject treated as a random effect. At each voxel, a one-sample t-test against a contrast  
400 value of zero gave us our estimate of statistical reliability. For whole brain analysis, we  
401 corrected for multiple comparison using cluster correction, with a cluster forming  
402 threshold of  $p < .001$  and an extent threshold calculated with SPM to set a family-wise  
403 error cluster level corrected threshold of  $p < .05$  (127 for PE>fixation; 267 for PE\*set size  
404 interaction). Note that these appropriately high cluster forming threshold ensures that  
405 parametric assumptions are valid and the rate of false positives are appropriate (Eklund  
406 et al. 2016; Flandin & Friston 2016).

## 407 **ROIs:**

408 Fronto-parietal network: As we did not have specific regional predictions regarding the  
409 WM component of learning, we defined broad fronto-parietal networks as ROIs that have  
410 been previously associated with a wide range of tasks involving cognitive control.  
411 Specifically, our first control network ROIs were defined by using left and right anterior

412 dorsal premotor cortex (prePMd: 8mm sphere around -38 10 34, (Badre & D'Esposito  
413 2007) as seeds in two separate “resting state” (task-free) seed-to-voxel correlation  
414 analyses in the CONN toolbox (<https://www.nitrc.org/projects/conn/>), and using the  
415 corresponding whole-brain connectivity to left and right prePMd, as our control network  
416 ROI. In order to confirm the robustness of our findings, we then ran a larger  
417 frontoparietal network ROI defined from a functionally neutral group (Yeo et al. 2011),  
418 along with a functionally defined ROI of the multiple demands network from (Fedorenko  
419 et al. 2013). All three of these frontoparietal ROIs yielded similar outcomes, thus  
420 confirming the robustness of our findings. We report here the results from the Yeo et al  
421 ROI as the widest, most neutral ROI.

422  
423 The striatum ROI was defined based on univariate activity for prediction error ( $p < .001$ ,  
424 uncorrected), masked by AAL definitions for putamen, caudate, and nucleus accumbens  
425 (Marsbar AAL structural ROIs: <http://marsbar.sourceforge.net/download.html>). We note  
426 that this ROI definition would be biased for assessing the effect of RPE in the striatum.  
427 However, this is not our goal as the relationship of RPE and striatum is established both  
428 in general from the prior literature, and in this study based on the corrected whole brain  
429 analysis (see Results). Rather, this ROI will be used to test the effects of set size and  
430 the interaction of set size with RPE, within regions maximally sensitive to RPE. As the  
431 set size variable is uncorrelated with that of RPE, this ROI definition does not bias either  
432 of these analyses.

433

434 For each ROI, a mean time course was extracted using the MarsBar toolbox  
435 (<http://marsbar.sourceforge.net/>). The GLM design was estimated against this mean time  
436 series, yielding parameter estimates (beta weights) for the entire ROI for each regressor  
437 in the design matrix.

438

439 **Whole brain Contrasts:** We focus on two main contrasts: 1) positive effect of RPE; 2)  
440 positive interaction of RPE and set size, to determine whether WM processes influence  
441 RPE signaling and whether such interactions relate to behavior. The first contrast is  
442 defined by considering the sum of the beta weights across all set sizes:  $\sum_{i=1:6} \beta_{PE(i)}$ ; we  
443 test whether this contrast value is significantly positive. The second contrast takes the  
444 linear contrast of the beta weights across set sizes by the set size:  $\sum_{i=1:6} (i-3.5) * \beta_{PE(i)}$  ;  
445 testing whether this contrast is positive signals a linear increase of RPE with set size.

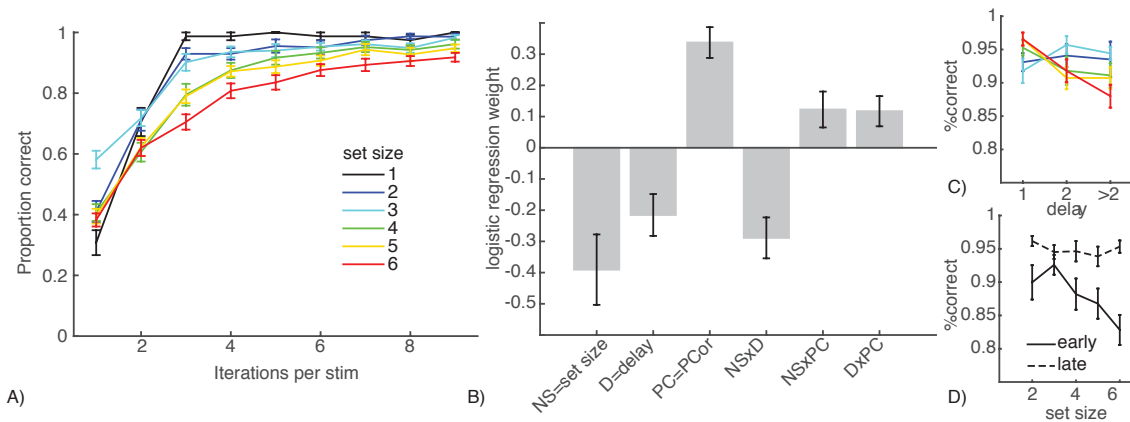
446 We also tested the opposite contrasts, as well as the linear effect of set size  $\sum_{i=1:6} (i-$   
 447  $3.5) * \beta_i$

448  
 449

450 **Interaction between set-size and RPE:** To investigate individual differences in the  
 451 interaction between set size and RPE, we assessed ROI markers of this interaction. We  
 452 computed this in one of three ways, each reflecting different assumptions: (A) a linear  
 453 contrast of set-size on RPE regression weight; (B) a contrast of high set size (4-6) vs.  
 454 low set size (1-3) on RPE regression weights (in case of a step function for e.g. above  
 455 vs. below capacity sets), and (C) Spearman rho of RPE weights across set-sizes, which  
 456 does not require linearity and is less susceptible to outliers than linear regression.  
 457 Despite slightly different assumptions, all three measures are highly correlated (all  
 458 rhos > 0.8, p < 10<sup>-4</sup>) and yielded qualitatively similar results. Because we observe that  
 459 results neither show linear changes across set sizes, nor a step function, we report  
 460 results using the measure defined as option C.

461 **Results:**

462 **Behavior:**



463

464 **Figure 2: Behavioral Results.** A) Proportion of correct choices as a function of how many times  
 465 a specific stimulus was encountered (i.e., learning curves), for each set size. B) Logistic  
 466 regression on factors that contribute to accuracy for a given image, including set size (NS), delay  
 467 since last previous correct choice for a given image (D), PCor (number of previous correct  
 468 choices for that image), and their interactions. C) Illustration of the interaction between delay and  
 469 set size. D) Illustration of the interaction between set size and PCor – early indicates PCor < 4, late  
 470 indicate PCor > 6. Error bars indicate standard error of the mean.

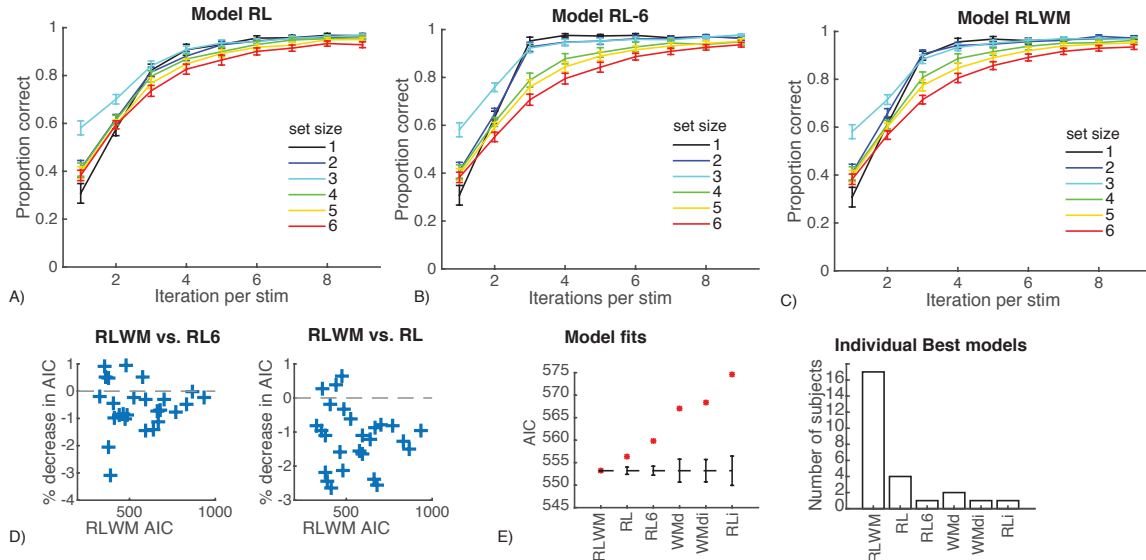
471

472 Behavioral results replicate our previous findings (Collins & Frank 2012; Collins et al.  
473 2014; Figure 2). Learning curves showed strong differences as a function of set size,  
474 despite the same number of encounters for each stimulus. Logistic regression analysis  
475 of subjects choices (Fig. 2B) showed main effects of reward history, delay and load,  
476 indicating that subjects were more likely to select the correct action with more previous  
477 correct experience for a given stimulus ( $t(25)=6.8$ ,  $p<10^{-4}$ ), and less likely to be correct  
478 with increasing set size ( $t(25)=-3.4$ ,  $p=.002$ ) and increasing delay (intervening trials since  
479 their last correct choice on this stimulus) ( $t(25)=-3.2$ ,  $p=.004$ ). There were also  
480 interactions between all pairs of factors, such that the delay effect was stronger in high  
481 load ( $t(25)=-4.4$ ,  $p=.0002$ , Fig. 2C), and the effects of load and delay both decreased with  
482 more correct reward history ( $t_s>2.1$ ,  $p_s<.05$ , Fig. 2D). The latter interaction is expected  
483 given the RLWM model's prediction that behavior transitions from WM (which is more  
484 sensitive to delay and load) to RL as a function of learned reliability.

#### 485 **Model fitting:**

486 Model fitting also confirmed our previous findings, showing that a computational model  
487 including two modules (RL and WM) explained subjects' behavior better than variants of  
488 a model assuming a single RL or WM process. Specifically, RLWM provided a  
489 significantly better AIC than RL6 ( $t(25)=3.9$ ,  $p=0.001$ ) and RL ( $t(25)=-6.6$ ,  $p<10^{-4}$ ), and  
490 individual AICs favored RLWM for a significant number of subjects (21/26 for RL6, sign  
491 test  $p=0.002$ ; 23/26 for RL,  $p<10^{-4}$ ). Model simulations show that a simple RL model  
492 cannot capture the behavior as well as RLWM or RL6, but note that RL6 needs too  
493 many parameters to appropriately capture behavior. Pure working memory models  
494 assuming changes in decay with experience, or interference, also cannot capture  
495 behavior as well as our hybrid RLWM model (Fig. 3E)

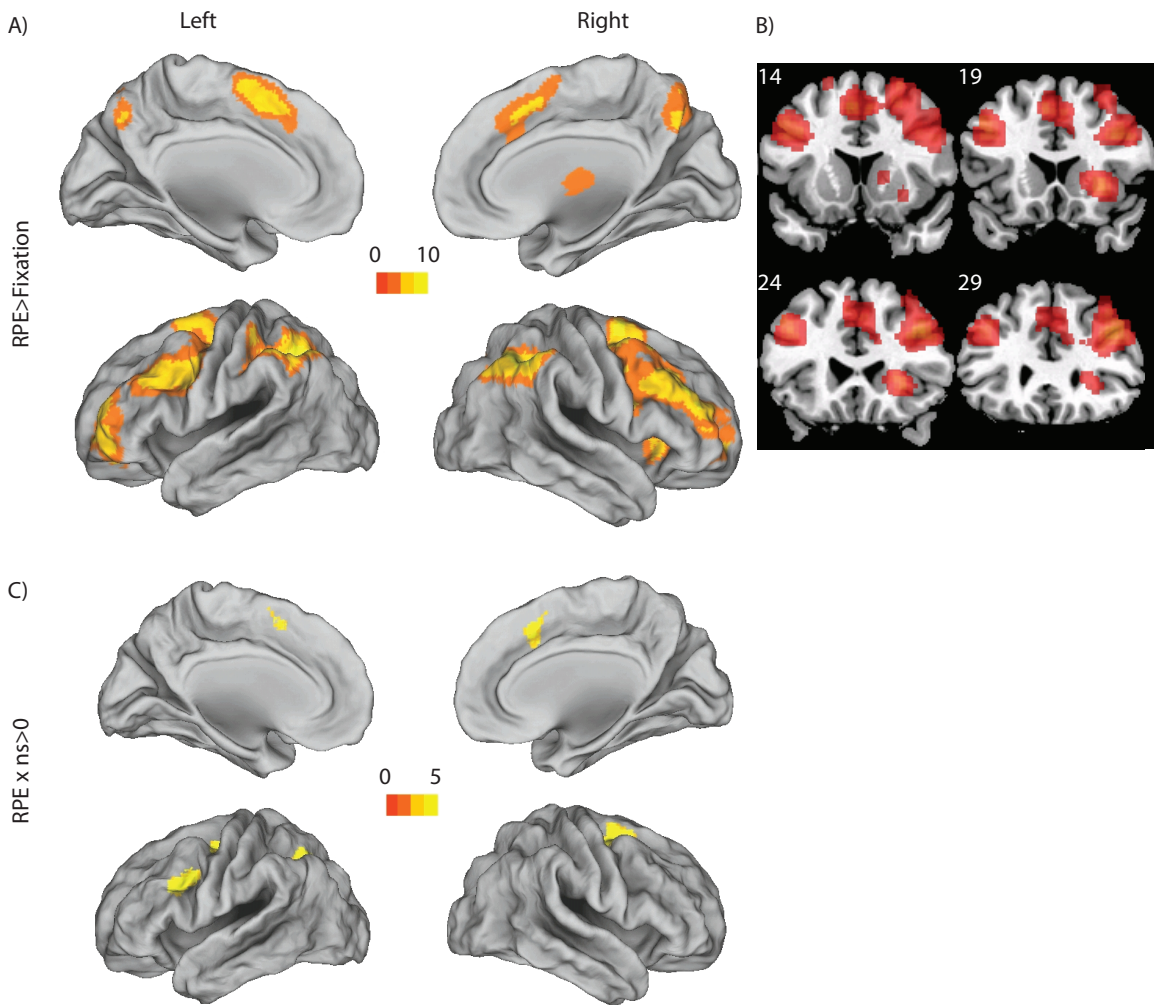




496  
 497  
 498  
 499  
 500  
 501  
 502  
 503  
 504  
 505  
 506

**Figure 3: Model Validation.** A-C) Proportion of correct responses as a function of how many times a specific stimulus was encountered, for each set size, for simulation of different models with individually fit parameters. Models were simulated 100 times per subject then averaged within subjects to represent this subject's contribution. Error bars indicate standard error of the mean across subjects. A) simple RL model including decay and different sensitivity to gains/losses. B) Identical model to A, with learning rate varying per set size. C) Model incorporating both RL and WM. D) Model comparisons show a significantly lower AIC for RLWM than RL6 or RL, for a significant number of subjects. Each cross indicates a single subject. E) Model comparison to other potential models show best fit for RLWM (see methods for other model names).

507 **Imaging results:**



508

509 **Figure 4: Whole brain effects of RPE and RPE x ns.** A-B) Regions positively correlated with  
510 RPE ( $p < .05$  cluster corrected). C) Regions showing a positive interaction of RPE with set size.  
511

512 Whole brain analysis showed increasing activity with set size in bilateral precuneus and  
513 decreasing activity in a network including bilateral superior frontal gyrus, bilateral angular  
514 gyrus and bilateral supramarginal gyrus (table 3), confirming that the set size  
515 manipulation is effective at differentially engaging large brain networks.

516

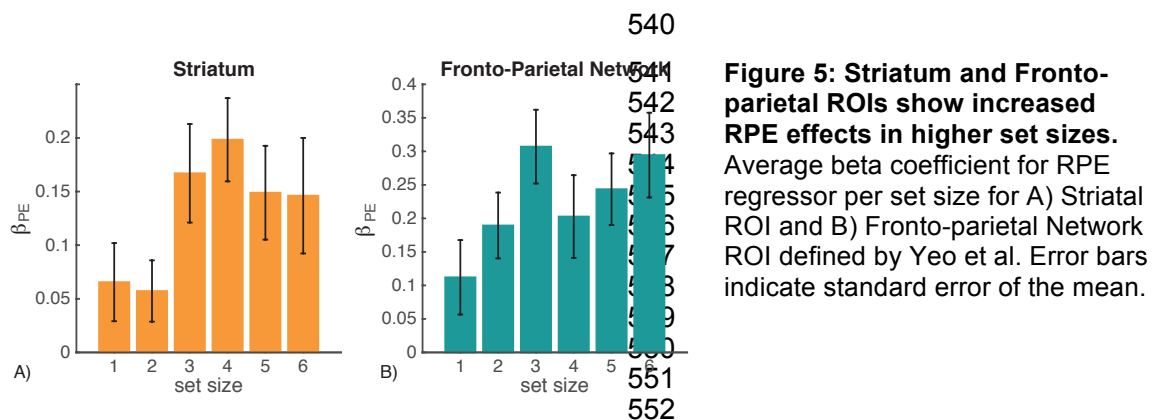
517 Whole brain analysis showed a distributed network that positively correlated with the  
518 parametric reward prediction errors (RPE) regressor. We verified RPE-related activation  
519 in the right caudate nucleus and thalamus (See table 1 for full results, figure 4B), as  
520 expected from the literature. Notably, the RPE network also includes regions of bilateral  
521 prefrontal and parietal cortex commonly observed in cognitive control tasks.

522

523 We next tested whether the RPE signal was homogeneous across set sizes in striatum,  
 524 as implicitly expected if striatal RL is independent of WM. To the contrary, we found a  
 525 significant positive interaction of set size with RPE ( $t(20)=2.4$ ,  $p=0.026$ ; figure 5B) in the  
 526 striatal ROI (see Methods). Note that this interaction reflects a *stronger* effect of RPE on  
 527 the striatal BOLD signal at higher set-sizes (i.e., under more cognitive load). This finding  
 528 supports the hypothesis that WM interacts with RL, showing blunted RL signals in low  
 529 set sizes (i.e., within the capacity of WM).

530

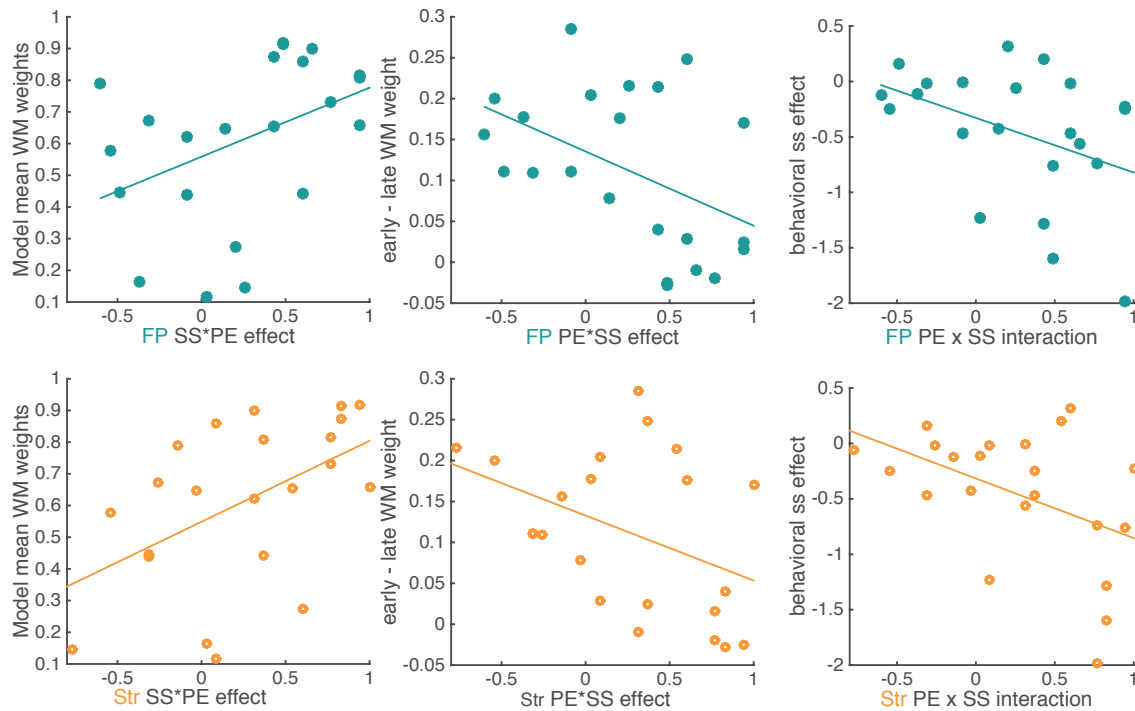
531 Next, we investigated whether other brain regions showed the same modulation of RPE  
 532 signaling by WM load. Whole brain analysis showed a positive linear interaction of set  
 533 size with RPE in left lateral prefrontal cortex and parietal cortex (MNI coordinates -38,  
 534 20, 28; table 2). Further investigation within an independent fronto-parietal network ROI  
 535 (Yeo et al. 2011) showed both a strong main effect of prediction error ( $t(20)=6.9$ ,  $p<10^{-4}$ )  
 536 and a significant interaction of set-size with RPE in the fronto-parietal ROI  
 537 ( $t(20)=2.3$ ,  $p=0.03$ ), a pattern similar to the striatum ROI. Again, RPE signaling was larger  
 538 with more WM load, possibly reflective of a common neuromodulatory signal in striatum  
 539 and cortex influenced by cognitive demands.



553 **Link to behavior:**

554 We hypothesized that the weaker RPE signals observed in low set-sizes might reflect an  
 555 interaction between WM and RL systems. Specifically, this may reflect the greater use of  
 556 WM, instead of RL, at low set sizes. This strategy could be because low set sizes do not  
 557 require RPE signaling: the most recent stimulus-action-outcome can be accessed from  
 558 memory. Thus, we predicted that those subjects relying more on WM would exhibit a  
 559 stronger neural interaction effect (i.e., they would show less homogeneity in their RPE  
 560 signals across set-sizes). To index WM contributions to choice, we use the

561 computational model-inferred weight of the WM module, averaged over all trials. Indeed,  
 562 we found that greater WM contributions to choices was significantly related to the set-  
 563 size effect on RPE signaling, both in striatum ( $\rho=0.55$ ,  $p=.01$ ), and the fronto-parietal  
 564 ROI ( $\rho=0.49$ ;  $p=0.02$ ; figure 6 left). Moreover, subjects who continued to rely on WM  
 565 with experience (i.e., exhibiting less transition to RL) also showed greater set-size  
 566 effects on RPE signaling in FP ( $\rho=-0.46$ ,  $p=0.03$ ) and marginally, in striatum ( $\rho=-$   
 567  $0.41$ ;  $p=0.06$ ; figure 6 middle). This may be due to the fact that for participants with  
 568 higher overall reliance on WM, the WM module is more reliable, and thus WM use  
 569 decreases less over learning. Indeed, the two indexes were negatively correlated ( $\rho=-$   
 570  $0.69$ ,  $p<10^{-3}$ ). The results were partly accounted for by differences in model fit capacity  
 571 parameter: subjects with higher capacity showed significantly stronger nsxRPE  
 572 interaction in FP ( $\rho=0.46$ ,  $p=0.03$ ), and marginally so in striatum ( $\rho=.41$ ,  $p=0.06$ );  
 573 Finally, we confirmed this effect was independent of the fit of the RLWM model by using  
 574 the logistic regression, and specifically the effect of set-size on accuracy (note that this  
 575 measure was, as expected, related to the one obtained by the computational model:  
 576 Spearman  $\rho=-.42$ ,  $p=0.05$ ). Indeed, the effect of set-size on accuracy was marginally  
 577 related to the set-size by RPE interaction in striatum ( $\rho=-0.4$ ,  $p=0.06$ , figure 6 right) and  
 578 FP ( $\rho=-0.41$ ,  $p=0.06$ ). Again, neural interactions were stronger for those subjects  
 579 exhibiting a stronger negative effect of set-size on behavior.



580

581 **Figure 6: Effect of set size on RPE in the fMRI signal is related to individual differences in**  
582 **behavior.** Left: average model-inferred mixture weight assigned to working memory over RL  
583 (“Model mean WM weight”) is significantly related to a stronger effect of set size in fronto-parietal  
584 ROI ( $\rho=0.49$ ,  $p=0.02$ ) and in the striatum ( $\rho=0.55$ ,  $p=0.01$ ). Middle: decrease in working memory  
585 weight from early (first 3 iterations) to late in a learning block (last 3 iterations) is significantly  
586 related to fMRI effect in FP ROI ( $\rho=-0.46$ ,  $p=0.03$ ), and marginally so in striatum ( $\rho=-0.41$ ,  
587  $p=0.06$ ). Right: the behavioral set-size effect is measured as the logistic regression weight of the  
588 set-size predictor; stronger behavioral effect is marginally related to a stronger neural effect in FP  
589 ROI ( $\rho=-0.41$ ,  $p=0.059$ ) and in striatum ROI ( $\rho=0.4$ ,  $p=0.063$ ).

## 590 Discussion:

591 We combined computational modeling and fMRI to investigate the contributions of two  
592 distinct processes to human learning: reinforcement learning and working memory. We  
593 replicated our previous results (Collins & Frank 2012; Collins et al. 2014) showing that  
594 these jointly play a role in decisions: computational models assuming a single learning  
595 process (either WM or RL) could not capture behavior adequately. We also replicated  
596 the widespread observation that the striatum and lateral prefrontal cortex are sensitive to  
597 reward prediction errors, a marker of RL. We made the novel observation that RL and  
598 WM are not independent processes, with the most commonly studied RL signal blunted  
599 under low WM load. Further, we found that the degree of interaction was related to  
600 individual differences in subjects’ use of WM: the more robustly subjects used WM for  
601 learning, the more they showed WM effects on RL signals.

602

603 The process of model-free reinforcement learning, as both a class of machine learning  
604 algorithms and as the neural network function implemented via dopamine-dependent  
605 plasticity in cortico-basal ganglia networks, is characterized by integration of rewards  
606 over time to estimate the value of different options, and a value dependent policy. Our  
607 behavioral results replicate our previous work showing that even in simple instrumental  
608 learning, we cannot account for human learning based only on the integrated history of  
609 reward. Instead, the influences of load and delay/intervening trials show that working  
610 memory also contributes to learning. That this influence decreases with experience  
611 supports a model where RL and WM modules are dynamically weighted according to  
612 their success in predicting observed outcomes.

613

614 We used computational modeling to disentangle the contributions of RL and WM to  
615 learning and to assess neural indicators of their interactions. We extracted the reward  
616 prediction error signal from the RL module, and confirmed in a model-based whole brain

617 fMRI analysis that striatum was sensitive to prediction errors, as established from a large  
618 literature (Pessiglione et al. 2006; Schönberg et al. 2007), as was a large bilateral fronto-  
619 parietal region (Daw et al. 2011). However, we found in both regions that sensitivity to  
620 RPE was modulated by set size, the number of items that subjects learned about in a  
621 given block. Specifically, the RPE signal was weaker in lower set sizes, in which  
622 subjects' learning was closest to optimal, and thus likely to mostly use WM. Thus, as  
623 noted in our earlier studies (Collins & Frank 2012; Collins et al. 2014), WM contributions  
624 to learning can confound measures typically attributed to RL. While the previous findings  
625 were limited to behavioral, genetic and computational model parameters, here we report  
626 for the first time that even neural RPE signals are influenced by WM. These results also  
627 imply that in other studies that do not manipulate WM load during learning, the  
628 contribution of WM to learning may yield inflated or blunted estimates of the pure RL  
629 process.

630

631 We further found that individual differences in the degree to which set-size modulated  
632 RPE signals correlated with the degree to which subjects relied on WM in their  
633 behavioral learning curves. Specifically, subjects with more robust use of WM showed  
634 more reliably blunted RPE signals in lower set sizes, supporting the interpretation that  
635 WM use induces weaker RPEs in the RL system. Further supporting this interpretation,  
636 we observed that subjects who continued to use WM with learning (i.e., showing less  
637 transition to RL) exhibited larger effects of set-size on RPE signaling.

638

639 One might expect to observe more reliable indicators of neural computations with easier  
640 tasks; our findings show the opposite. These results thus strongly hint at a mechanism  
641 by which WM and RL interact beyond the competition for control of action (Poldrack et  
642 al. 2001), and specifically at a mechanism by which WM interferes with RL  
643 computations. How might this interference occur? One possibility is that the two  
644 processes compete not only for guiding action, but also more generally, for example  
645 based on their reliability in a given environment. Such interference would mean that in  
646 conditions in which WM performs better than RL (eg. early in learning for low set sizes),  
647 WM inhibits the whole RL mechanism and thus weakens its characteristic neural signals,  
648 such as RPEs. Another possible explanation for the observed interference is cooperative  
649 interaction, where WM modifies the reward expectations in the RL system. This would  
650 lead – when WM was working well – to higher expectations than would be computed by

651 pure RL, and thus to weaker RPEs. Future research will need to distinguish these  
652 possibilities. There may be other interpretations of the change in RPE signaling with set  
653 size, besides our interpretation as an interaction between the RL and WM processes.  
654 However, given that behavioral fits strongly implicate separate WM and RL processes in  
655 learning (see above and previous studies), and that WM is sensitive to load in other  
656 paradigms with similar profiles, this remains the most parsimonious explanation. Note  
657 that this interaction also makes other behavioral predictions suggesting that  
658 reinforcement value learning is actually enhanced under high WM load; we have recently  
659 confirmed this prediction using a novel task building on this line of work (Collins et al,  
660 submitted).

661

662 Our results are related to recent work on sequential decision making and learning that  
663 highlighted the role of a model-free module (similar to our RL model), and of a model-  
664 based module, responsible for representing stimulus-action-outcome transitions and  
665 using them to plan decisions (Doll et al. 2015). This latter module has been linked to  
666 cognitive control and is weakened under load (Otto et al. 2013), suggesting that it may  
667 require WM. Moreover, both WM use in the current task and model-based processing in  
668 the sequential task are related to the same genetic variant associated with prefrontal  
669 catecholaminergic function (Collins & Frank 2012; Doll et al. 2016). Notably, (Daw et al.  
670 2011) showed that RPEs in the striatum were modulated by model-based values, a  
671 result that may support our collaborative hypothesis. However, we demonstrate such  
672 interaction even in paradigms that are traditionally thought to involve purely “model-free”  
673 RL. As there is no sequential dependence between trials, learning in our paradigm does  
674 not require learning a transition model or planning. Indeed, we could adequately capture  
675 learning curves for individual set-sizes using a purely model-free RL model (Collins &  
676 Frank 2012; Collins et al. 2014), with decreasing learning rates across set sizes, but this  
677 model has more parameters than RLWM and cannot capture the nuanced effects of e.g.,  
678 delay and set-size interactions. Thus, our results show that learning in very simple  
679 environments that appear to require purely model-free learning still recruits executive  
680 functions, with working memory contributing to learning and interfering with the putative  
681 dopaminergic RL process. Our results show a similar pattern of RPE activations for  
682 subcortical and lateral prefrontal cortex areas, a common finding in published studies  
683 (e.g. Badre & Frank 2012; Frank & Badre 2012), possibly reflecting a common  
684 dopaminergic input to both regions (Bjorklund & Dunnett 2007).

685

686 We investigated the role of working memory using set-size as a proxy. However, this  
687 leaves open some questions and may limit some of our interpretations. In particular, set-  
688 size affects the overall load of working memory, but is also predictive of higher delays  
689 between repetitions of the same stimulus. While our analyses tease apart load from  
690 delay, the delay itself comprises both a temporal component (number of seconds over  
691 which working memory could decay passively), and a discrete component (number of  
692 intervening trials that may interfere with working memory). Our paradigm did not  
693 manipulate those two factors to make them maximally decorrelated, and cannot  
694 distinguish their relative contributions to the effect of delay on behavior. Furthermore, by  
695 focusing on set size as the marker of WM, we cannot distinguish between a “tonic”, or  
696 slowly tuned interference of WM in RL computation, vs. a more “phasic”, trial-by-trial  
697 adjustment of their role and interaction between them. A target for future research is  
698 increasing the experimental paradigm’s capacity to carefully disentangle delay from load,  
699 allowing us to better understand the dynamics of interactions between RL and WM.

700

701 We focused on WM as an alternative learning mechanism from RL, with an a priori  
702 interest in regions of the cognitive control network in lateral frontal and parietal cortices.  
703 However, regions involved in long-term memory (LTM), such as medial temporal lobe  
704 (MTL) and hippocampus, could also play an important role: rote memorization of explicit  
705 rules is in the prime domain of LTM, others have shown trade-offs for learning between  
706 LTM and striatal based learning (Poldrack et al. 2001), and WM itself is often difficult to  
707 distinguish from LTM (Ranganath & Blumenfeld 2005; D’Esposito & Postle 2015). Our  
708 results are consistent with LTM having a role in learning: indeed, we observe a negative  
709 correlation between RPE and activation in a network of regions including MTL (table 2),  
710 indicating higher activation early in learning (Poldrack et al. 2001). However,  
711 computational modeling shows that the second learning component we extract is  
712 capacity limited, supporting our interpretation of this component as mainly WM.  
713 Nevertheless, future research is needed to more carefully dissociate the role of WM from  
714 LTM in reinforcement learning.

715

716 Learning is a key factor in humans improving their abilities, skills, and fitting to our  
717 quickly changing environments. Understanding what distinct cognitive and neurological  
718 components contribute to learning is thus essential, in particular to study differences in



719 learning across individuals. Many neurological and psychiatric disorders include learning  
720 impairments (Huys et al. 2016). To precisely understand how learning is affected by  
721 these conditions, we must be able to reliably extract separable cognitive factors,  
722 understand how these factors interact, and link them to their underlying neural  
723 mechanisms. Our results provide a first step toward clarifying how we trade off working  
724 memory and integrative value learning to make decisions in simple learning  
725 environments, and how these processes may interfere with each other.

726 Badre, D. & D'Esposito, M., 2007. Functional magnetic resonance imaging evidence for a hierarchical  
727 organization of the prefrontal cortex. *Journal of cognitive neuroscience*, 19(12), pp.2082–99.  
728 Badre, D. & Frank, M.J., 2011. Mechanisms of Hierarchical Reinforcement Learning in Cortico-Striatal  
729 Circuits 2: Evidence from fMRI. *Cerebral cortex (New York, N.Y. : 1991)*, pp.1–10.  
730 Bjorklund, A. & Dunnett, S.B., 2007. Dopamine neuron systems in the brain: an update. *Trends in*  
731 *Neurosciences*, 30(5), pp.194–202.  
732 Botvinick, M.M., Niv, Y. & Barto, A.C., 2009. Hierarchically organized behavior and its neural  
733 foundations: a reinforcement learning perspective. *Cognition*, 113(3), pp.262–80.  
734 Burnham, K.P. & Anderson, D.R., 2002. *Model Selection and Multi-Model Inference: A Practical*  
735 *Information-Theoretic Approach (Google eBook)*, Springer.  
736 Collins, a. G.E. et al., 2014. Working Memory Contributions to Reinforcement Learning Impairments in  
737 Schizophrenia. *Journal of Neuroscience*, 34(41), pp.13747–13756.  
738 Collins, A. & Koechlin, E., 2012. Reasoning, Learning, and Creativity: Frontal Lobe Function and Human  
739 Decision-Making J. P. O'Doherty, ed. *PLoS Biology*, 10(3), p.e1001293.  
740 Collins, A.G.E. & Frank, M.J., 2013. Cognitive control over learning: Creating, clustering, and  
741 generalizing task-set structure. *Psychological Review*, 120(1), pp.190–229.  
742 Collins, A.G.E. & Frank, M.J., 2012. How much of reinforcement learning is working memory, not  
743 reinforcement learning? A behavioral, computational, and neurogenetic analysis. *The European*  
744 *journal of neuroscience*, 35(7), pp.1024–35.  
745 Collins, A.G.E. & Frank, M.J., Opponent actor learning (OpAL): Modeling interactive effects of striatal  
746 dopamine on reinforcement learning and choice incentive.  
747 D'Esposito, M. & Postle, B.R., 2015. The Cognitive Neuroscience of Working Memory. *Annual Review of*  
748 *Psychology*, 66(1), pp.115–142.  
749 Daw, N.D. et al., 2011. Model-based influences on humans' choices and striatal prediction errors. *Neuron*,  
750 69(6), pp.1204–15.  
751 Daw, N.D. & Doya, K., 2006. The computational neurobiology of learning and reward. *Current opinion in*  
752 *neurobiology*, 16(2), pp.199–204.  
753 Doll, B.B. et al., 2015. Model-based choices involve prospective neural activity. *Nature neuroscience*,  
754 (February), pp.1–9.  
755 Doll, B.B. et al., 2016. Variability in Dopamine Genes Dissociates Model-Based and Model-Free  
756 Reinforcement Learning. *Journal of Neuroscience*, 36(4), pp.1211–1222.  
757 Eklund, A., Nichols, T.E. & Knutsson, H., 2016. Cluster failure: Why fMRI inferences for spatial extent  
758 have inflated false-positive rates. *Proceedings of the National Academy of Sciences*, 113(28),  
759 pp.7900–7905.  
760 Fedorenko, E., Duncan, J. & Kanwisher, N., 2013. Broad domain generality in focal regions of frontal and  
761 parietal cortex. *Proceedings of the National Academy of Sciences of the United States of America*,  
762 110(41), pp.16616–21.  
763 Flandin, G. & Friston, K.J., 2016. Analysis of family-wise error rates in statistical parametric mapping  
764 using random field theory.  
765 Frank, M.J., Seeberger, L.C. & O'reilly, R.C., 2004. By carrot or by stick: cognitive reinforcement learning  
766 in parkinsonism. *Science (New York, N.Y.)*, 306(5703), pp.1940–3.  
767 Guitart-Masip, M. et al., 2012. Go and no-go learning in reward and punishment: interactions between  
768 affect and effect. *NeuroImage*, 62(1), pp.154–66.  
769 Hart, a. S. et al., 2014. Phasic Dopamine Release in the Rat Nucleus Accumbens Symmetrically Encodes a  
770 Reward Prediction Error Term. *Journal of Neuroscience*, 34(3), pp.698–704.  
771 Huys, Q.J.M., Maia, T. V & Frank, M.J., 2016. Computational psychiatry as a bridge from neuroscience to  
772 clinical applications. *Nature Neuroscience*, 19(3), pp.404–413.  
773 Jocham, G., Klein, T. a & Ullsperger, M., 2011. Dopamine-mediated reinforcement learning signals in the  
774 striatum and ventromedial prefrontal cortex underlie value-based choices. *The Journal of*  
775 *neuroscience : the official journal of the Society for Neuroscience*, 31(5), pp.1606–13.  
776 Montague, P.R., Dayan, P. & Sejnowski, T.J., 1996. A framework for mesencephalic dopamine systems  
777 based on predictive Hebbian learning. *The Journal of neuroscience : the official journal of the Society*  
778 *for Neuroscience*, 16(5), pp.1936–47.  
779 Nassar, M.R. & Frank, M.J., 2016. Taming the beast: Extracting generalizable knowledge from  
780 computational models of cognition. *Current Opinion in Behavioral Sciences*, 11, pp.49–54.  
781 Otto, a R. et al., 2013. Working-memory capacity protects model-based learning from stress. *Proceedings*

782           *of the National Academy of Sciences of the United States of America*, 110(52), pp.20941–6.  
783 Pessiglione, M. et al., 2006. Dopamine-dependent prediction errors underpin reward-seeking behaviour in  
784 humans. *Nature*, 442(7106), pp.1042–5.  
785 Petrides, M., 1985. Deficits on conditional associative-learning tasks after frontal- and temporal-lobe  
786 lesions in man. *Neuropsychologia*, 23(5), pp.601–614.  
787 Poldrack, R. a et al., 2001. Interactive memory systems in the human brain. *Nature*, 414(November),  
788 pp.546–550.  
789 Ranganath, C. & Blumenfeld, R.S., 2005. Doubts about double dissociations between short- and long-term  
790 memory. , 9(8).  
791 Schönberg, T. et al., 2007. Reinforcement learning signals in the human striatum distinguish learners from  
792 nonlearners during reward-based decision making. *The Journal of neuroscience : the official journal*  
793 *of the Society for Neuroscience*, 27(47), pp.12860–7.  
794 Schultz, W., 2002. Getting formal with dopamine and reward. *Neuron*, 36(2), pp.241–63.  
795 Schwarz, G., 1978. Estimating the Dimension of a Model. *The Annals of Statistics*, 6(2), pp.461–464.  
796 Sutton, R.S. & Barto, A.G., 1998. *reinforcement learning*, MIT Press.  
797 Yeo, B.T.T. et al., 2011. The organization of the human cerebral cortex estimated by intrinsic functional  
798 connectivity. *Journal of neurophysiology*, 106, pp.1125–1165.  
799  
800

## Figure Legends

**Figure 1: Experimental Protocol.** At the beginning of each block, subjects were shown for 10 s the set of stimuli they would see in that block. In this example, Block 1 uses color patches for stimuli and has a set size  $n_s = 2$ ; Block n uses shapes and has  $n_s = 6$ . Each trial included the presentation of a stimulus for 0.5s followed by a blue fixation cross until subject pressed one of three buttons, or up to 1.5s after trial onset. Button press caused the fixation cross to turn white. Feedback was presented for 1s, and came 1.5s after trial onset. Feedback consisted of the words correct or incorrect in green and red, respectively. The inter-trial interval consisted of a white fixation cross with jittered duration to allow trial by trial event-related analysis of fMRI signal. Blocks set sizes varied between 1 and 6, and the order was randomized across subjects.

**Figure 2: Behavioral Results.** A) Proportion of correct choices as a function of how many times a specific stimulus was encountered (i.e., learning curves), for each set size. B) Logistic regression on factors that contribute to accuracy for a given image, including set size (NS), delay since last previous correct choice for a given image (D), PCor (number of previous correct choices for that image), and their interactions. C) Illustration of the interaction between delay and set size. D) Illustration of the interaction between set size and PCor – early indicates  $PCor < 4$ , late indicate  $PCor > 6$ . Error bars indicate standard error of the mean.

**Figure 3: Model Validation.** A-C) Proportion of correct responses as a function of how many times a specific stimulus was encountered, for each set size, for simulation of different models with individually fit parameters. Models were simulated a 100 times per subject then averaged within subjects to represent this subject's contribution. Error bars indicate standard error of the mean across subjects. A) simple RL model including decay and different sensitivity to gains/losses. B) Identical model to A, with learning rate varying per set size. C) Model incorporating both RL and WM. D) Model comparisons show a significantly lower AIC for RLWM than RL6 or RL, for a significant number of subjects. Each cross indicates a single subject. E) Model comparison to other potential models show best fit for RLWM (see methods for other model names).

**Figure 4: Whole brain effects of RPE and RPEXns.** A-B) Regions positively correlated with RPE ( $p < .05$  cluster corrected). C) Regions showing a positive interaction of RPE with set size.

835 **Figure 5: Striatum and Fronto-parietal ROIs show increased RPE effects in higher set**  
836 **sizes.** Average beta coefficient for RPE regressor per set size for A) Striatum ROI and B) Fronto-  
837 parietal Network ROI defined by Yeo et al. Error bars indicate standard error of the mean.

838

839 **Figure 6: Effect of set size on RPE in the fMRI signal is related to individual differences in**  
840 **behavior: Effect of set size on RPE in the fMRI signal is related to individual differences in**  
841 **behavior.** Left: average model-inferred mixture weight assigned to working memory over RL  
842 (“Model mean WM weight”) is significantly related to a stronger effect of set size in fronto-parietal  
843 ROI ( $r=0.49$ ,  $p=0.02$ ) and in the striatum ( $r=0.55$ ,  $p=0.01$ ). Middle: decrease in working memory  
844 weight from early (first 3 iterations) to late in a learning block (last 3 iterations) is significantly  
845 related to fMRI effect in FP ROI ( $r=-0.46$ ,  $p=0.03$ ), and marginally so in striatum ( $r=-0.41$ ,  $p=0.06$ ).  
846 Right: the behavioral set-size effect is measured as the logistic regression weight of the set-size  
847 predictor; stronger behavioral effect is marginally related to a stronger neural effect in FP ROI ( $r=-$   
848  $0.41$ ,  $p=0.059$ ) and in striatum ROI ( $r=0.4$ ,  $p=0.063$ ).

849  
850

## fMRI activations from Prediction Error contrasts

851 **Table 1: Main effect of RPE**

852 *All clusters reliable at  $p < .05$  corrected. Coordinates are the center of mass in MNI.*

853 **A) Contrast: Main Effect of RPE > Fixation**

Region	BA	Extent (voxels)	x	y	z	Peak t-val
Right Angular Gyrus	7	3202	34	-60	42	10.83
	40		46	-52	44	9.2
Right Inferior Parietal Gyrus	40		42	-42	40	10.21
Left Superior Parietal Gyrus	7	3317	-30	-54	44	10.43
Left Angular Gyrus	40		-46	-48	56	10.32
Left Inferior Parietal Gyrus	40		-42	-42	42	9.45
Right Superior Frontal Sulcus	6	12409	20	2	62	9.64
Right Middle Frontal Gyrus	46		38	36	30	8.84
Left Superior Frontal Gyrus	6		-24	-6	62	8.15
Left Middle Frontal Gyrus	11	1686	-30	56	4	7.78
Left Lateral Orbital Gyrus	46		-40	56	-2	6.96
Left Anterior Orbital Gyrus	11		-24	44	-14	6.42
Right Putamen		955	28	22	0	6.99
Right Thalamus			12	-10	10	5.15
Right Pallidum			12	0	6	4.36
Right Precuneus	7	731	6	-64	40	6.32
	7		8	-66	58	5.21

854 **B) Contrast: Main Effect of RPE < Fixation**

Region	BA	Extent (voxels)	x	y	z	Peak t-val
Right Superior Occipital Gyrus	18	9715	16	-92	24	10.22
Left Superior Occipital Gyrus	18		-16	-96	18	8.73
Right Inferior Lingual Gyrus	30		-10	-48	-6	8.9
Left Cingulate Gyrus (subgenual)	11	2264	-4	28	-12	8.52
	25		-2	18	-8	7.34

Left Superior Frontal Gyrus	10		-8	58	2	7.24
Left Middle Temporal Gyrus	20	2543	-56	-8	-18	6.69
Left Supramarginal Gyrus	48		-36	-36	22	6.26
Left Superior Temporal Gyrus	38		-34	8	-20	6.24
Right Precentral Sulcus	4	1248	26	-30	66	6.21
Right Postcentral Gyrus	4		36	-26	72	6.13
Right Precentral Gyrus	4		52	-12	58	5.62
Right Superior Temporal Gyrus	38	336	30	10	-28	6.08
Right Middle Temporal Gyrus	21		50	2	-26	5.56
	21		58	0	-24	4.76
Right Cingulate Gyrus	23	516	6	-20	44	5.64
Right Superior Frontal Gyrus	6		12	-18	62	5.03
Right Cingulate Sulcus	4		10	-16	54	4.56
Right Superior Temporal Gyrus	48	935	54	-4	4	5.6
Right Lateral Fissure	48		50	4	-6	5.11
Right Lateral Fissure/Insular Gyrus	48		40	-14	20	5.04

855

856

857 **Table 2: set-size \* RPE interaction**

**Contrast: RPE Parametric Increasing With Set Size**

Region	BA	Extent (voxels)	x	y	z	Peak t-val
Left Superior Precentral Sulcus	44	725	-46	10	36	5.69
Left Inferior Frontal Sulcus	48		-38	20	28	5.16
Left Middle Frontal Gyrus	6		-32	2	38	4.57
Right Superior Frontal Gyrus	6	689	18	4	54	5.42
	32		6	22	46	4.3
Left Superior Frontal Gyrus	6		-6	10	50	4.08
Left Intraparietal Sulcus	7	463	-26	-66	44	5.28
	7		-30	-58	46	5.24
	19		-26	-68	34	4.59

858

859

860 **Table 3: Main effect of set size**

861 **A) Contrast: Set Size Parametric Increasing**

Region	BA	Extent (voxels)	x	y	z	Peak t-val
Left Precuneus	7	1948	-6	-72	44	6.98
Left Angular Gyrus	40		-32	-50	36	6.4
Right Precuneus	7		12	-70	44	5.22

862

863 **B) Contrast: Set Size Parametric Decreasing**

Region	BA	Extent (voxels)	x	y	z	Peak t-val
Right Superior Frontal Gyrus	9	1344	14	58	34	6.64
Left Superior Frontal Gyrus	9		-12	46	42	4.69
	10		-4	58	28	4.5
Left Supramarginal Gyrus	40	447	-64	-44	34	6.31
Left Angular Gyrus	39		-52	-70	28	5.83
	40		-60	-52	40	4.52
Right Angular Gyrus	40	255	58	-52	44	5.41
	22		62	-54	28	4.41
Right Supramarginal Gyrus	40		64	-46	36	5.19
Right Superior Frontal Gyrus	8	239	14	20	62	5.15
	9		10	38	52	4.42
Left Superior Frontal Sulcus	8	255	-24	22	58	4.9
Left Middle Frontal Gyrus	46		-24	18	40	4

864

865



866

A)	K	$\alpha$	$\phi_{WM}$	$\rho$	$\phi_{RL}$	$\epsilon$	pers
Mean	4.08	0.07	0.29	0.86	0.05	0.03	0.34
(std)	(0.98)	(0.13)	(0.31)	(0.18)	(0.05)	(0.03)	(0.31)
Median	4	0.03	0.18	0.94	0.05	0.03	0.25
Min - max	2-5	0.01- 0.5	0-1	0.42-1	0-0.21	0-0.14	0.02-1

867

868

B)	K	$\alpha$	$\phi_{WM}$	$\rho$	$\phi_{RL}$	$\epsilon$
$\alpha$	ns					
$\phi_{WM}$	ns	0.77				
$\rho$	ns	-0.65	-0.77			
$\phi_{RL}$	ns	0.83	0.69	-0.62		
$\epsilon$	ns	ns	ns	ns	ns	
pers	ns	ns	ns	ns	ns	ns

869

870

871 Table 4) RLWM model fit parameters. A) Parameter statistics B) Correlation between  
 872 parameters. ns indicates non-significant correlation ( $p < .05$ , corrected for multiple  
 873 comparisons.

874