

1 Proactive and reactive construction of memory- 2 based preferences

3 **Jonathan Nicholas^{1,2,3}, Nathaniel D. Daw^{4,5}, and Daphna Shohamy^{1,2,6}**

4 ¹Department of Psychology, Columbia University, New York, NY, USA

5 ²Mortimer B. Zuckerman Mind, Brain, Behavior Institute, Columbia University, New York, NY, USA

6 ³Department of Psychology, New York University, New York, NY, USA

7 ⁴Department of Psychology, Princeton University, Princeton, NJ, USA

8 ⁵Princeton Neuroscience Institute, Princeton University, Princeton, NJ, USA

9 ⁶The Kavli Institute for Brain Science, Columbia University, New York, NY, USA

10 **Abstract**

11 We are often faced with decisions we have never encountered before, requiring us to infer
12 possible outcomes before making a choice. Computational theories suggest that one way to make
13 these types of decisions is by accessing and linking related experiences stored in memory. Past
14 work has shown that such memory-based preference construction can occur at a number of
15 different timepoints relative to the moment a decision is made. Some studies have found that
16 memories are integrated at the time a decision is faced (reactively) while others found that
17 memory integration happens earlier, when memories are encoded (proactively). Here we offer a
18 resolution to this inconsistency. We demonstrate behavioral and neural evidence for both
19 strategies and for how they tradeoff rationally depending on the associative structure of memory.
20 Using fMRI to decode patterns of brain responses unique to categories of images in memory, we
21 found that proactive memory access is more common and allows more efficient inference.
22 However, participants also use reactive access when choice options are linked to more numerous
23 memory associations. Together, these results indicate that the brain judiciously conducts
24 proactive inference by accessing memories ahead of time in conditions when this strategy is most
25 favorable.

26 **Introduction**

27 Some decisions are made repeatedly, offering the opportunity to learn directly about an option's
28 value through past experiences with its outcome. However, decisions often consist of a choice
29 between options whose outcomes have not been directly experienced before. Computational
30 theories of planning suggest that one way to approach such decisions is by knitting together
31 separate relevant memories through mental simulation¹⁻³. The ability to flexibly combine
32 information in this way is central to intelligence: it frees us from having to decide based on direct
33 trial-and-error experience alone and enables us to make inferences and to plan novel courses of
34 action using cognitive maps or internal models⁴⁻⁸.

35 The process of drawing inferences requires accessing relevant memories and recombining or
36 integrating across them to build new relationships. Studying memory access is therefore one way
37 to shed light on the covert mechanisms that give rise to inferential choice. Yet previous work
38 attempting to probe this connection has left open a critical gap in our understanding of how and
39 when memory integration supports inference. In particular, some studies have claimed that
40 memories are accessed at the time a choice is faced^{2,9,10}, while other studies have found that
41 memory access occurs much earlier, when relevant memories are first encoded^{11,12}. These two
42 approaches differ not just in the timepoint of memory access, but also point to distinct
43 mechanisms. Specifically, integrating memories during a decision requires “on the fly” processing,
44 which is likely to take time, whereas integrating memories earlier suggests that the new model for
45 inference already exists when a choice is later made, yielding more efficient decisions^{11,13,14}. It
46 has been suggested, but not yet empirically tested, that there may be some normative explanation
47 for the variation between these two approaches¹⁵. In the present study, we aimed to address this
48 gap by studying both possibilities in a single experimental design. We sought to first confirm the
49 normative advantages that early memory access confers and then to investigate how changing
50 the structure of memory access can rationally shift this process to happen later, at decision time.

51 The role of memory integration in inference is often studied with multi-phase tasks that first seed
52 relevant associative memories and then test whether people integrate them when probed to make
53 decisions. A classic task in this vein, which we build upon here, is *sensory preconditioning*¹⁶. In
54 sensory preconditioning, participants are first trained to associate two stimuli that occur in

55 succession (A→B). Then, in a separate phase, the B stimulus is associated with reward. The
56 critical question is whether people infer that the A stimulus is also associated with reward. This is
57 tested in the final decision phase, when participants are asked to choose between A and another
58 control stimulus (which is equally familiar but lacks the indirect reward association). Humans and
59 non-human animals alike tend to prefer A despite never directly experiencing its association with
60 reward^{11,12,14,16}. Neural studies of sensory preconditioning and similar tasks have revealed two
61 potential mechanisms, each predicting memory integration either before or during choice, that
62 may lead to this same behavioral effect.

63 The most typical explanation for inference in tasks like sensory preconditioning, widely assumed
64 in theories of decision making dating back to Tolman⁸, envisions that choosing A reflects
65 prospective mental simulation at decision time: in this case, retrieving the B-reward association
66 when evaluating whether to choose A. This, in turn, is thought to be a minimal case of our more
67 general capacity for constructive, deliberative forward planning, embodied in theories of model
68 based reinforcement learning, which iteratively evaluate candidate actions prospectively over
69 multiple steps using a learned internal model of task contingencies. By examining neural
70 signatures of memory retrieval, it has been possible to investigate how memory access actually
71 relates to successful model-based inference. Yet, these studies have yielded mixed support for
72 this account. Some evidence suggests that both humans and non-human animals engage in
73 prospective retrieval at decision time, and that this pattern is associated with inferential
74 performance^{4,9,10,17–19}. However, there is also evidence that associative recall may occur long
75 before a decision is ever faced^{11,12,20–23}.

76 These latter findings imply a second explanation for inference in these tasks: that the value of
77 options may be pre-computed when relevant information like reward is first encoded, thereby
78 preempting the need for evaluating potential outcomes later at choice time. In some studies of
79 sensory preconditioning, for instance, it has been found that when B is presented during reward
80 learning, A is concurrently retrieved and directly associated with reward^{11,12}. Such a strategy is
81 feasible because, at this time, participants have already been provided with all of the components
82 necessary to form a complete model of the task. Perhaps analogously, in rodent spatial navigation
83 tasks, hippocampal place cells often briefly represent trajectories in front of the animal^{17–19}, a
84 potential substrate for prospective evaluation. However, otherwise similar “replay” events can
85 instead reflect backward or altogether nonlocal trajectories at the time of reward^{24–27}, potentially
86 supporting a spatial analogue of the alternative inference strategy.

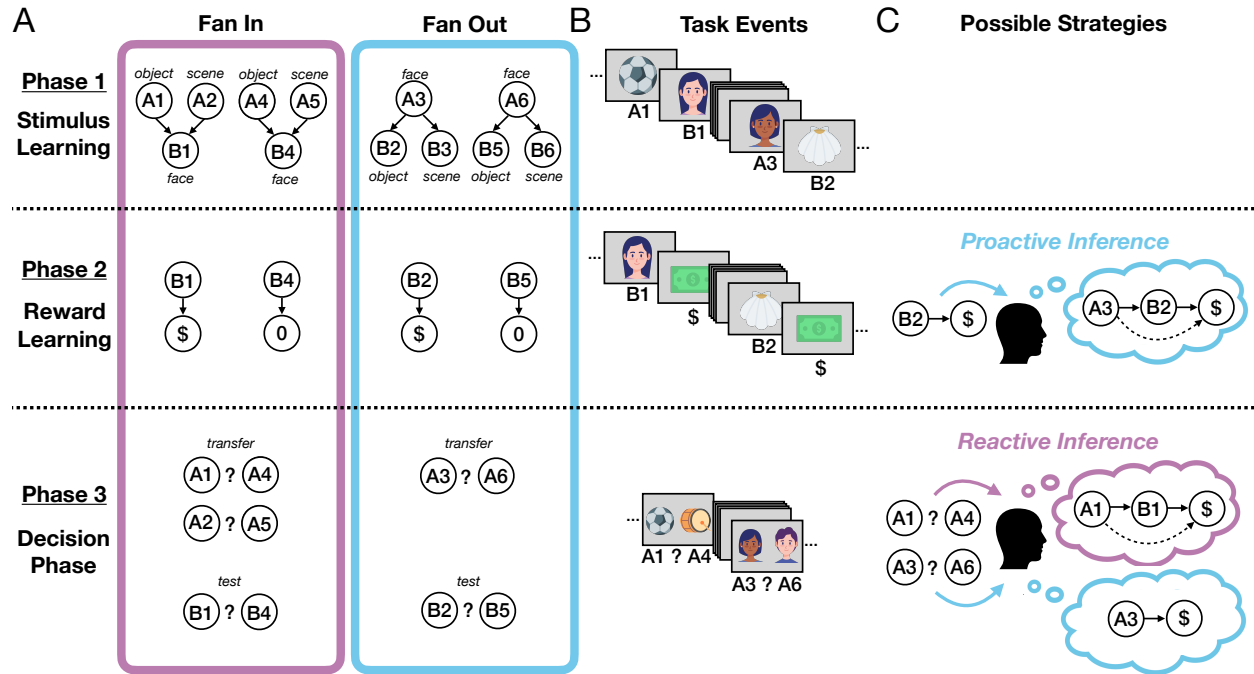
87 An emerging idea is that these different inference mechanisms may be special cases of a more
88 general set of computations that share the common goal of integrating memories to infer action
89 values, but that access memories at different times: either *proactively* before they are needed or
90 *reactively*, once required for choice^{15,28}. This in turn raises questions about how these strategies
91 are balanced or adaptively deployed, and whether such control might explain variable results
92 across studies. Indeed, the possibility of proactive computation implies that the brain must
93 somehow be judicious about which memories it accesses, and when, since there are so many
94 possible later actions that might be contemplated.

95 This idea, while compelling, is still largely untested, and raises a number of questions about how
96 and when different strategies are deployed, which we aimed to address in this study. First, is it
97 indeed the case that a proactive memory access strategy exists and can support inferential choice
98 equivalent to a reactive one? Second, what are the tradeoffs of the two approaches: if access
99 occurs proactively, does it indeed reduce the need for computation at decision time? Finally, do
100 people rely differentially on this strategy at times when it would be sensible to do so?

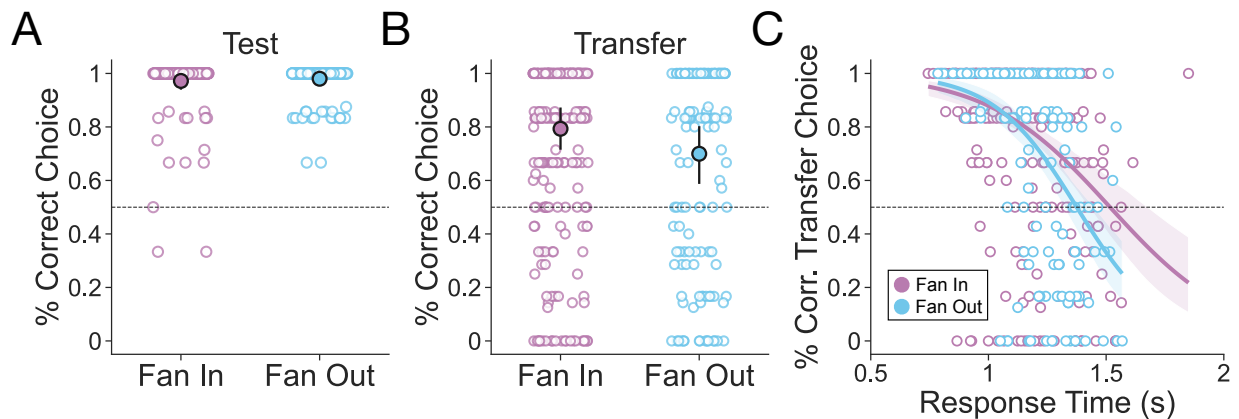
101 Here we aimed to answer these questions by attempting to alter participants' reliance on proactive
102 inference. We had three primary hypotheses. First, we aimed to confirm earlier (but inconsistently
103 reported) results that sensory preconditioning can be solved with proactive memory access at the
104 time of reward learning. Next, because proactive inference offers the advantage of a pre-
105 computed value association, we hypothesized that this approach may allow for more efficient
106 future decisions—i.e. decisions that are faster and more accurate. Finally, we hypothesized that
107 reliance on this strategy would adapt under different circumstances, which we operationalized by
108 manipulating how difficult it is to access and integrate relevant memories. Drawing upon a rich
109 tradition of research on associative memory²⁹, we reasoned that having multiple relevant
110 associations with an experience should, at any timepoint, induce competition between them,
111 making their retrieval for use in inference less likely.

112 To test these hypotheses, we developed a novel learning and decision making task based on
113 sensory preconditioning, and measured memory retrieval at multiple timepoints of this task while
114 scanning participants with fMRI (**Figure 1**). To capture putative reactivation of associations in
115 memory in the service of inference, we exploited the fact that viewing different visual categories
116 (e.g. faces, scenes, and objects) elicits unique activity in visual cortex^{10,11,30,31}. We used images
117 from these different categories for each of the different stimuli, which allowed us to measure
118 whether reactivation of associated images in memory occurred during either reward learning,
119 signifying proactive inference, or during decision making, signifying reactive inference. We
120 predicted that proactive memory access during reward learning should result in more efficient
121 later choices, and that reactive memory access during choice itself should have the opposite
122 effect.

123 To address our third hypothesis specifically, we further varied the number of competing
124 associations with a given stimulus by training participants on stimulus-stimulus relationships
125 under two different conditions. In one condition, two *antecedent* stimuli each predicted a single
126 *consequent* stimulus; we refer to this as the *Fan In* condition. By contrast, in the *Fan Out* condition,
127 a single *antecedent* predicted two possible *consequents*. The logic of this manipulation is that the
128 *Fan In* condition induces greater retrieval competition between memories of antecedent stimuli
129 when the consequent stimulus is presented during the reward learning phase. We therefore
130 predicted that there should be increased reliance on reactive inference for stimuli in the *Fan In*
131 condition relative to *Fan Out* condition.



132
 133 **Figure 1. Task design and inference strategies. A) Task structure.** Participants (n=39) underwent fMRI
 134 scanning while completing a three-part experiment with two different conditions, based on sensory
 135 preconditioning. The phases were similar for both conditions, which differed only in their specific associative
 136 structure. In phase one, *stimulus learning*, participants learned associations between several pairs of
 137 images (faces, scenes, or objects). Unknown to participants, there were two types of trials governing how
 138 these associations appeared. *Fan In* trials consisted of one of two possible antecedent A images followed
 139 by one consequent B image. *Fan Out* trials consisted of one antecedent A image followed by one of two
 140 possible consequent B images. Example categories for each image are shown here, and this was
 141 counterbalanced across participants. In phase two, *reward learning*, participants learned that a subset of
 142 consequent B images led to a reward, while others did not lead to reward. Finally, in phase three, the
 143 *decision phase*, participants chose between two images. Choices between consequent B images were
 144 used as *test* trials, whereas choices between antecedent A images were used as *transfer* trials. **B) Example**
 145 **events.** An example of the sequence of task events seen by participants in each phase. **C) Possible**
 146 **inference strategies.** Participants can engage in either of two inference strategies: proactive inference, at
 147 the time of reward learning, or reactive inference, at the time of the decision. During decision making,
 148 proactive inference does not require the integration of a memory with value, as this association has already
 149 been performed during reward learning. Due to differences in the number of competing antecedent
 150 memories at reward learning, we expected reactive inference to be used more for Fan In stimuli.



151
152 **Figure 2. Participants successfully learned and transferred across both conditions, but the**
153 **relationship between speed and accuracy differed across conditions. A)** Test decisions (i.e. those
154 between images that were directly associated with reward or neutral outcomes during reward learning)
155 were highly accurate, reflecting successful learning for both conditions. **B)** Transfer decisions (i.e. those
156 between images that were indirectly associated with reward or neutral outcomes via the stimulus learning
157 phase) were also highly accurate, indicating successful inference for both conditions. Filled points represent
158 group-level means whereas white points represent means for each pair of images seen by participants.
159 Error bars are 95% confidence intervals. **C)** The relationship between the proportion of accurate transfer
160 choices and reaction time for each image pair revealed that faster decisions were more accurate and that
161 this relationship was stronger for the Fan Out condition, in which the structure was more amenable to
162 proactive integration. Lines represent regression fits and bands represent 95% confidence intervals. All
163 visualizations show data at the stimuli level, and statistical analyses were conducted using mixed effects
164 models that additionally assessed these effects within each participant while accounting for variation across
165 participants.

166 Results

167 Behavioral evidence for proactive inference and its modulation by retrieval 168 competition

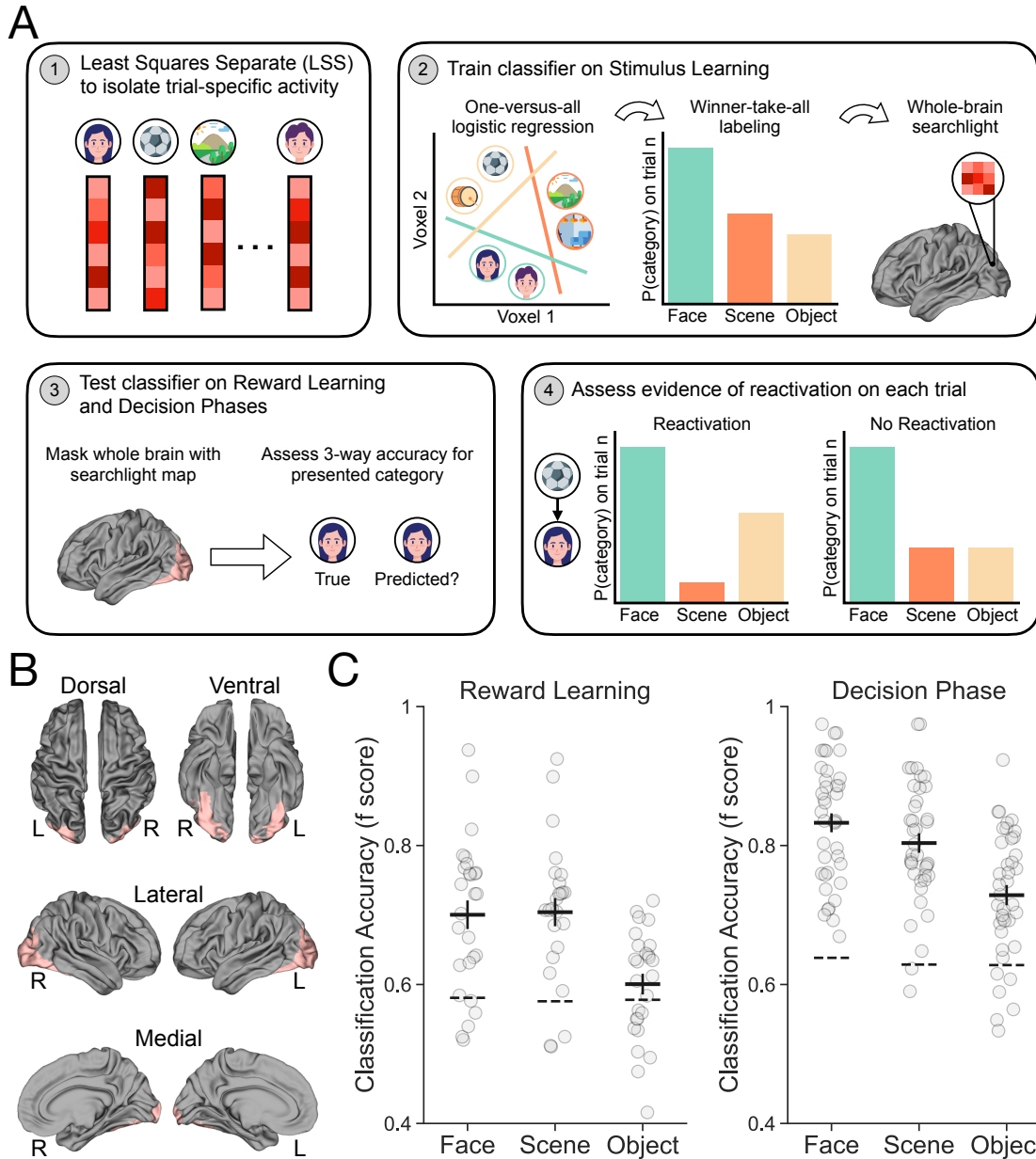
169 We first examined whether participants learned to directly associate consequent stimuli with
170 reward, and whether they transferred value to associated antecedent images. To assess this, we
171 analyzed participants' *test* and *transfer* choices during the decision phase. On test choices,
172 participants chose between consequent images that were directly associated with either a reward
173 or neutral outcome during the reward learning phase. Participants were highly accurate and
174 tended to choose the rewarded consequent image over the neutral consequent image ($\beta_0 =$
175 5.009 , $95\% CI = [4.085, 6.279]$; **Figure 2A**). There was no difference between the Fan In and
176 Fan Out conditions ($\beta_{condition} = 0.321$, $95\% CI = [-1.251, 2.128]$), indicating that participants
177 learned similarly in both.

178 Next, we examined participants' transfer choices during the decision phase (**Figure 2B**). These
179 decisions consisted of choosing between antecedent images that were paired with consequent
180 images during the initial stimulus learning phase. Critically, successful transfer of value to these
181 images involves relying on memory for the paired association. We found that participants tended
182 to choose the antecedent image that was paired with the rewarded consequent image ($\beta_0 =$
183 2.075 , $95\% CI = [1.283, 2.896]$), indicating that most participants used memory to transfer value.
184 There was no difference in transfer performance between Fan In and Fan Out choices
185 ($\beta_{condition} = 0.572$, $95\% CI = [-0.157, 1.284]$), demonstrating that the manipulation of

186 associative structure between conditions had no effect on the degree to which value was
187 transferred.

188 Having established that participants infer the value of associated antecedent images in both
189 conditions, we next sought to gain initial insights into *when* memories are accessed to support
190 this value transfer. We aimed to differentiate between two possible strategies for inference, each
191 occurring at different timepoints in our task: either proactively at reward learning or reactively at
192 decision time. One hypothetical hallmark of proactive inference is that it should promote accuracy
193 without further integration at choice time, resulting in faster transfer decisions. Thus, if its
194 deployment varies across stimuli, it predicts an unusual inverted speed-accuracy relationship
195 whereby faster decisions tend also to be more accurate. In contrast, successful reactive inference
196 requires integration at choice time, resulting in slower transfer decisions and (to the extent its
197 deployment governs successful performance) a more typical relationship between slower
198 decisions and higher accuracy.

199 Overall, we found that choices reflecting memory-based transfer were faster ($\beta_{rt} =$
200 -0.611 , 95% $CI = [-0.945, -0.287]$; **Figure 2C**), suggesting that participants tended to infer
201 proactively. Yet we also found that this relationship was stronger in the Fan Out than the Fan In
202 condition ($\beta_{condition:rt} = -0.465$, 95% $CI = [-0.937, -0.017]$), consistent with our expectation
203 that the Fan In condition is less amenable to proactive inference. Together, these behavioral
204 findings suggest that while proactive inference dominated performance overall, reactive inference
205 may have been more commonly observed in the Fan In than the Fan Out condition.



206

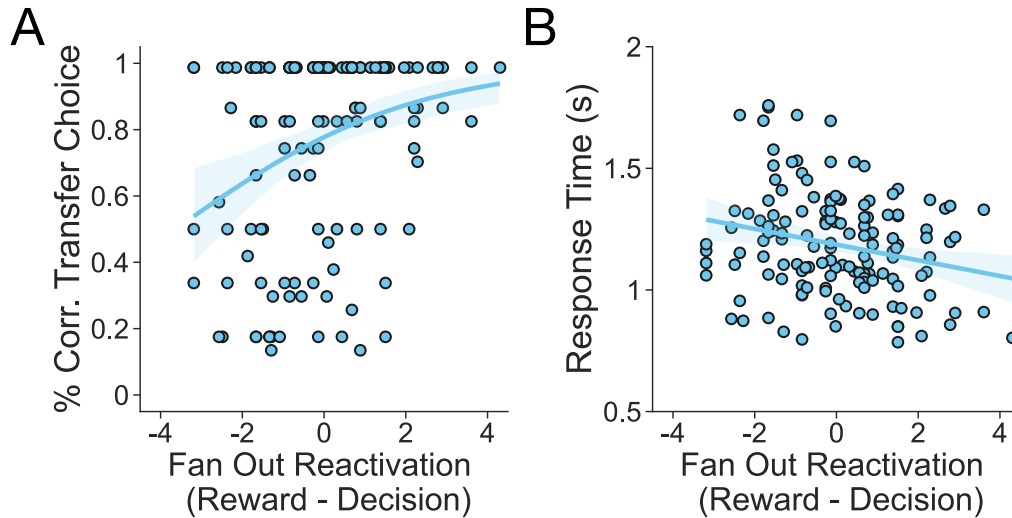
207 **Figure 3. Multivariate pattern analysis methodology and decoding accuracy.** **A)** MVPA analyses
 208 consisted of four primary steps. Step 1: Least Squares Separate³² was used to isolate a beta map for each
 209 trial and participant across all phases of the experiment. These betas were then used as input for the MVPA
 210 pipeline. Step 2: A searchlight analysis consisting of a one versus all three-way logistic regression was then
 211 used to identify voxels that could discriminate between all three categories during the stimulus learning
 212 phase. Step 3: Voxels identified during the previous step were then used to mask the whole brain during
 213 testing of the classifier on the reward learning and decision phases. Step 4: Evidence of reactivation on
 214 each trial was then assessed by ranking the individual category probabilities accordingly. **B)** Group-level
 215 whole-brain maps (FDR corrected; $q < 0.05$) of voxels that discriminate between all three categories above
 216 chance. **C)** Classification accuracy for the decoding model trained on the stimulus learning phase and
 217 tested on the reward learning and decision phases. Accuracy is shown here as the weighted F-score. Points
 218 represent accuracy for each participant and the thick line represents group-level average accuracy. Dotted
 219 lines represent the 95th percentile of a permutation distribution over test category labels.

220 **Neural evidence for proactive and reactive inference and their modulation by** 221 **retrieval competition**

222 While examining participants' choices allowed us to assess the different behavioral signatures of
223 proactive and reactive inference, choice behavior alone cannot capture when exactly memories
224 were accessed throughout the task. To gain further insight into when memories were recalled to
225 support inference, we used fMRI to obtain a neural signature of memory reactivation at different
226 timepoints in our task (**Figure 3A**). We first used runs of fMRI data collected from the stimulus
227 learning phase to train a classifier to distinguish between each image category: faces, scenes or
228 objects. We then tested this classifier on activity from the reward learning and decision making
229 phases, and assessed its ability to identify the category of the image that was presented to
230 participants. As expected, voxels that differentiated accurately between categories were located
231 primarily across the bilateral occipito-temporal cortex (**Figure 3B**). When tested on the reward
232 learning and decision making phases, the classifier accurately differentiated each category from
233 the others (Faces: $\beta_0 = 0.161$, 95% $CI = [0.134, 0.189]$; Scenes: $\beta_0 = 0.151$, 95% $CI =$
234 $[0.123, 0.180]$; Objects: $\beta_0 = 0.066$, 95% $CI = [0.041, 0.093]$; **Figure 3C**).

235 With a classifier in hand that could distinguish between each category based on BOLD activity
236 patterns seen during the reward learning and decision phases, we were poised to assess the
237 degree to which memories were reactivated for inference, and when. Specifically, to measure
238 memory reactivation, we examined the individual category probabilities from the classifier on
239 every trial, and identified those in which the probability of the *associated* image category (as
240 opposed to the presented category) was particularly high (see **Methods**). This analysis allowed
241 us to label every trial as one in which reactivation of the relevant associated category in memory
242 was either likely or unlikely.

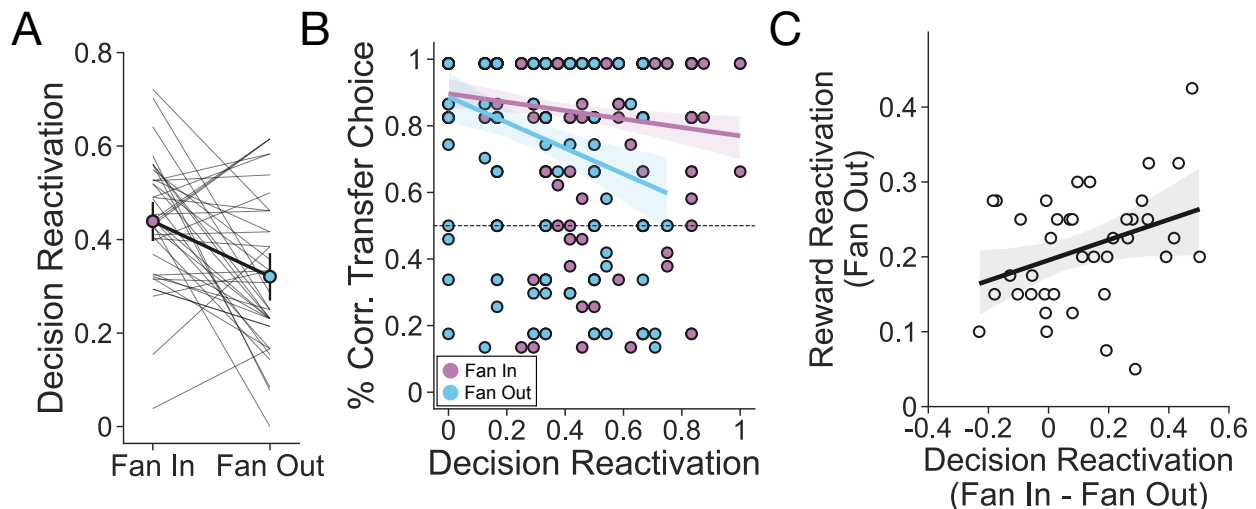
243 To determine whether memories were accessed in accordance with the patterns of inference we
244 observed behaviorally, we focused on three main goals for the analyses. First, because
245 participants' choice behavior at transfer suggested a tradeoff between speed and accuracy most
246 consistent with proactive inference, we sought to examine whether greater memory reactivation
247 during the reward learning phase indeed results in more efficient (faster and more accurate)
248 choices. Second, because we found that this effect was weaker during Fan In compared to Fan
249 Out transfer choices (when there was relatively more retrieval competition between memories
250 during reward learning and less during decision making), we sought to determine whether this
251 behavioral shift was supported by different memory access patterns across conditions. Third, we
252 predicted that it would be most strategic for participants to proactively infer prior to choice time for
253 Fan Out trials, but to reactively infer at choice time for Fan In trials and therefore tested this by
254 characterizing individual differences in memory access between participants.



255
256 **Figure 4. Proactive inference improves decision making ability.** Greater memory reactivation at reward
257 time relative to decision time - a marker of proactive inference - is associated with more effective transfer
258 decisions. **A)** Correct transfer decisions were more likely for pairs with greater memory reactivation during
259 reward learning relative to decision making. **B)** Response times were marginally faster for pairs with greater
260 memory reactivation during reward learning relative to decision making. Points represent average
261 performance for each image pair seen by participants. Lines represent regression fits and bands represent
262 95% confidence intervals. Visualizations show data at the stimuli level, and statistical analyses were
263 conducted using mixed effects models that additionally assessed these effects within each participant while
264 accounting for variation across participants.

265 To first examine whether memory access during reward learning leads to more efficient choices,
266 we quantified the difference in memory reactivation during image viewing at reward learning and
267 decision time. This yielded an index of proactive inference for each pair of images. We focused
268 on the Fan Out condition because the design allowed us to measure reactivation for this condition
269 at both of these time points (for the Fan In condition, the design only allows measuring reactivation
270 at decision time; see **Methods**). When there was more evidence of proactive inference - i.e. when
271 memory reactivation was greater at the time of reward learning relative to that of decision making
272 - transfer choices were both more accurate ($\beta_{\Delta reactivation} = 0.302$, 95% CI = [0.0384, 0.593]) and
273 marginally faster ($\beta_{\Delta reactivation} = -37.902$, 90% CI = [-75.273, -2.508], 95% CI =
274 [-82.823, 3.180]; **Figure 4**). This result suggests that using memory to transfer value via proactive
275 inference offers the advantage of more efficient choices in the future.

276 We next examined whether the Fan In and Fan Out conditions affected memory access patterns,
277 focusing on the time of choice because this was the timepoint at which we were able to assess
278 reactivation in both conditions (see **Methods**). In line with participants' behavior, we found that
279 during the decision phase, memories of associated consequent images were reactivated more
280 frequently for Fan In than Fan Out transfer decisions ($\beta_{condition} = 0.119$, 95% CI =
281 [0.051, 0.184]; **Figure 5A**). This result indicates that reactive inference is more likely to occur
282 when proactive inference is disadvantaged due to increased competition between memories for
283 retrieval prior to choice.



284
285 **Figure 5. Reactive inference is more likely in the Fan In than Fan Out condition.** **A)** Reactivation during
286 the decision phase was greater for Fan In than Fan Out trials. Filled points represent group-level means,
287 error bars are 95% confidence intervals, and thin lines represent individual subject slopes. **B)** Greater
288 memory reactivation at decision time, a marker of reactive inference, is associated with less effective
289 transfer decisions for Fan Out but not Fan In image pairs. Points represent average performance for each
290 image pair seen by participants. Lines represent regression fits and bands represent 95% confidence
291 intervals. **C)** Participants who showed greater reactivation for Fan In relative to Fan Out trials during
292 decision making also preferentially reactivated more for Fan Out trials during reward learning. Points
293 represent individual subjects, the line represents a linear regression fit, and the band represents a 95%
294 confidence interval.

295 The behavioral findings showed that reactive inference was associated with a lower proportion of
296 successful transfer decisions in the Fan Out relative to the Fan In condition (**Figure 2C**). This
297 effect may reflect the fact that, due to competition, proactive inference is easier and reactive
298 inference is correspondingly harder, making it less likely to be successful in the Fan Out condition.
299 We therefore predicted that the neural measure of memory reactivation at decision time should
300 likewise be associated with less successful value transfer in the Fan Out condition. Indeed, we
301 found that Fan Out transfer decisions were less accurate when antecedent memories were
302 reactivated at decision time ($\beta_{\text{reactivation}} = -0.300$, 95% CI = [-0.625, -0.001]; **Figure 5B**).
303 Further, no such effect was found in the Fan In condition ($\beta_{\text{reactivation}} = -0.086$, 95% CI =
304 [-0.255, 0.075]; **Supplementary Figure 1**). This result lends additional support to the
305 interpretation that the manipulation of associative structure increased participants' relative use of
306 reactive inference in the Fan In condition.

307 Finally, we assessed the idea that it would be strategic to proactively infer prior to choice time for
308 Fan Out trials, and to reactively infer at choice time for Fan In trials. We examined whether
309 individuals who tend to reactivate memories more for Fan In relative to Fan Out trials at decision
310 time also reactivated memories more for Fan Out trials during the reward learning phase. That is,
311 we asked whether participants' ability to appropriately deploy one of these strategies also
312 predicted appropriate deployment of the other. We found that this was indeed the case—
313 participants who reactivated memories more for Fan In transfer decisions relative to Fan Out
314 transfer decisions also reactivated memories for Fan Out stimuli at reward learning
315 ($\beta_{\Delta\text{reactivation}} = 0.027$, 95% CI = [0.003, 0.050]; **Figure 5C**). This result suggests that those
316 participants who were most sensitive to the presence of retrieval competition at either timepoint
317 strategically modulated when they accessed their memories to perform inference.

318 Discussion

319 Research on sequential decision making has found that the process of linking memories to
320 support inference is well described by theories based on model-based reinforcement learning^{4–}
321 ^{6,9,10}. Numerous studies have shown that memory-based inference can occur at a number of
322 different timepoints relative to the moment a decision is made^{10–12,18,19,22,23,33,34}. However, the
323 conditions that lead some memories to be accessed later than others have remained unclear.
324 Here we developed a task to directly test multiple hypotheses about the purpose and adaptability
325 of memory access in inference. Using fMRI to decode patterns of BOLD response unique to the
326 categories of images in memory, we found that participants primarily accessed memories
327 proactively, but this pattern was also sensitive to the situation: when a choice option had multiple
328 past associations, participants were more likely to defer inferring relationships between stimuli
329 and outcomes until decisions were made. We also found neural and behavioral evidence that
330 reinstating memories prior to decision making facilitates faster and more accurate inference,
331 suggesting that it is adaptive to plan in advance when possible. Together, these results indicate
332 that the brain judiciously conducts proactive inference, accessing memories proactively in
333 conditions when this is most favorable.

334 These findings add empirical support to predictions from computational work on model-based
335 reinforcement learning. This type of learning grants agents the ability to compose sequences of
336 simulated experience from a world model in order to discover the consequences of never
337 experienced actions. The process of simulating potential actions can occur in a forward manner,
338 by adding up expected immediate rewards over some future trajectory, or backwards, by
339 propagating value information from a destination state to a series of predecessors. These patterns
340 have been formalized by a number of different algorithms^{15,35,36}, and recent work has provided a
341 rational account of when each is most useful for decision making¹⁵.

342 Specifically, Mattar and Daw (2018) theorized that memories that are particularly likely to increase
343 future expected reward will be prioritized for reinstatement during inference and planning.
344 Formally, they proposed that the expected utility of accessing a past experience can be
345 decomposed into the product of two terms: need and gain. Need quantifies how likely an
346 experience is to be encountered again, and gain captures how much reward is expected from
347 improved decisions if that experience is reinstated. A critical feature of this model is that when the
348 need term dominates, memories tend to be accessed reactively at choice time, but if instead the
349 gain term dominates, memories tend to be accessed proactively following the receipt of reward.
350 The present findings generally support this theory. In particular, gain increases for an antecedent
351 when choices fan out, favoring proactive memory access, while need increases for consequents,
352 promoting reactive choice-time memory access, as they fan in. Thus, antecedents that are
353 associated with many consequents (i.e. that fan out) are more likely to be reinstated upon learning
354 that a consequent is rewarded, because there is much to gain from updating future decisions
355 made upon future encounters with the antecedent. Likewise, antecedents which deterministically
356 lead to a single consequent (i.e. that fan in) imply greater need for that consequent, and are more
357 likely to reinstate it at decision time. Importantly, while our findings are consistent with this
358 framework, they were also designed to be predicted by more intuitive, qualitative reasoning about
359 the degree of competition among different memories, and so go beyond any single theory of
360 prioritization for memory access.

361 In addition to findings from sensory preconditioning demonstrating that humans use memories for
362 inference of decision making, a number of other studies have shown that memory-based
363 inference may also take place offline, during periods of rest or sleep before choice. This approach

364 is advantageous because it offloads computation to otherwise unoccupied time. In humans, fMRI
365 research has revealed that memories are reactivated during periods of rest following reward^{20,21}
366 and that this reinstatement can enhance subsequent memory performance^{37,38}. Importantly, such
367 offline replay of past memories during rest has been shown to facilitate later integrative
368 decisions^{22,23}. Parallel work in rodents has demonstrated that hippocampal replay of previously
369 experienced spatial trajectories is observed during rest and sleep³⁹⁻⁴¹, and that rewarded
370 locations are replayed more frequently²⁵. These results indicate another way in which inferences
371 may be drawn offline, well before constituent memories are needed for choice. An important
372 direction for future work will be to see if rational considerations, such as sensitivity to competition
373 between memories, also affect the likelihood, or targets of, offline inference.

374 Another avenue for future study that we did not touch upon here involves the role of dopamine in
375 supporting the integration of memories with reward to guide behavior. Although the dopaminergic
376 system has traditionally been thought to support habitual learning from direct experience, recent
377 results suggest that dopamine may also support integrative evaluations of actions through the
378 flexible combination of past experience^{42,43}. Our task may provide an opportunity to further
379 elucidate the role of dopamine in this process. Despite being solved in different ways, both of the
380 conditions in our task are dependent upon the flexible expression of knowledge about stimulus
381 associations. Therefore, if dopamine is necessary for the acquisition of model-based associations,
382 as has been recently suggested⁴³, we expect it to be involved in both conditions equally. This
383 prediction could, for example, be tested by examining how integrative choice behavior in the
384 present task is affected by dopamine depletion in Parkinson's disease.

385 Recent behavioral work in humans has also shown that a strategy for backwards prediction similar
386 to the proactive inference strategy we measured here provides benefits for a number of different
387 types of decisions⁴⁴. In particular, this study demonstrated that such a strategy is relied upon more
388 often in environments where the number of states that follow a starting state outnumber those
389 that precede a rewarded state. Using a similar manipulation coupled with more direct assays of
390 strategy use, our results provide convergent evidence for this idea. Our study further enhances
391 understanding of proactive and reactive approaches to inference by grounding each of these
392 strategies in the mechanisms of memory.

393 Separately, one shortcoming of our study was that, due to our design, we were unable to isolate
394 memory reactivation when consequent images from the Fan In condition were presented during
395 reward learning. In practice, this limited our contrasts between conditions to decision time; and
396 our contrasts between timepoints to the Fan Out condition. This was because our metric of
397 memory reactivation was conservative in the sense of being selective to the specific relevant
398 candidate for classification. In particular, in addition to the category actually present on the screen
399 being most strongly decoded, we required that the relevant associate be more strongly activated
400 than the irrelevant foil to declare reactivation successful. However, at reward time in the Fan In
401 condition, both categories are relevant associates, so this comparison was not possible. One
402 possibility to skirt this issue in future work may be to present images of a fourth entirely unrelated
403 category. We did not pursue this direction in the present study to minimize the complexity of the
404 design. Future complementary work may explore these issues in more depth in order to allow for
405 cleaner measurement of reactivation when antecedent images fan in during reward learning.

406 In conclusion, we have demonstrated that the statistical structure of training experience impacts
407 whether inference from memory occurs before or during decision making. This finding suggests
408 that standard model-based prospective inference is not unique, but is instead one of a general
409 set of computations that access memory at different times. Together, these findings further help

410 to explain why different studies have observed memory integration to support choice at different
411 times, and suggest that different inference strategies may be recruited depending on their efficacy
412 for the task at hand.

413 **Materials and Methods**

414 **Participants**

415 A total of 40 participants (19 M, 21 F) between the ages of 18 - 35 were recruited from the
416 Columbia University community. Participants were right-handed, had normal or corrected-to-
417 normal vision, took no psychiatric medication, and had no diagnosis of psychological disorders.
418 One participant was removed from the analyses due to both failing to understand the instructions
419 of the task and missing responses on over half of the decision trials. The remaining 39 participants
420 had a mean age of 21.9 with a range of 19-35 and were included in the reported sample. Informed
421 consent was obtained at the beginning of the session and all experimental procedures were
422 approved by the Columbia University Institutional Review Board.

423 **Experimental Task**

424 Participants completed a three-part associative learning task while undergoing an fMRI scan. In
425 the first phase of the experiment, *stimulus learning*, participants were tasked with learning pairs
426 of images presented one at a time. Each trial consisted of a single image (A; 1.5s), followed by a
427 interstimulus interval in which a fixation cross was displayed (exponentially jittered with mean=3s,
428 min=0.5s, max=12s), followed by another image (B; 1.5s), and finally an intertrial interval in which
429 another fixation cross was displayed (exponentially jittered with mean=3s, min=0.5s, max=12s).
430 In order to ensure that participants were paying attention, they were asked to press a button box
431 with their index finger for the first image and with the middle finger for the second image in a pair.
432 Participants were shown 16 different pairs of images 5 times each for a total of 80 trials. Trials
433 were spread across two runs of 40 trials each. Images came from one of three categories, either
434 a face, a scene, or an object. In the second phase of the experiment, *reward learning*, participants
435 were tasked with learning that a subset of B images from the stimulus learning phase led
436 deterministically to reward, while another subset of images led deterministically to a neutral
437 outcome. Each trial consisted of a single image (1.5s), followed by an interstimulus interval in
438 which a fixation cross was displayed (2s), followed by the outcome (either a dollar bill or a gray
439 rectangle; 1.5s), and then finally an intertrial interval (exponentially jittered with mean=2.5s,
440 min=0.5s, max=10s). Participants were told to withhold a response for the image and to respond
441 with their index finger when a dollar was shown and with their middle finger when a gray rectangle
442 was shown. Participants saw each of 8 images 10 times for a total of 80 trials. Trials were spread
443 across two runs of 40 trials each. During the third and final phase of the experiment, the *decision*
444 *phase*, participants were tasked with deciding between two images of the same category (either
445 A v. A or B v. B) presented on the screen simultaneously. Each trial consisted of a choice
446 (max=2s), a confirmation in which a green rectangle appeared around their choice (2s-reaction
447 time), and then an intertrial interval (exponentially jittered with mean=2.5s, min=0.5s, max=10s).
448 Participants pressed with their index finger to choose the image on the left hand side of the screen
449 and with their middle finger to choose the image on the right hand side of the screen. Participants
450 made 78 choices across a single run of this phase. Interstimulus intervals and trial ordering was
451 optimized to minimize the correlation between events throughout each phase of the task.

452 The pairs of stimuli presented throughout the experiment fell into one of two conditions that were
453 unknown to participants: Fan Out and Fan In trials. Fan Out trials consisted of one A image that
454 could be followed by either of two B images, while Fan In trials consisted of either of two A images

455 followed by one B image. During stimulus learning, eight pairs of images fanned in, while another
456 eight fanned out. Of the eight pairs from each condition, there were two pairs of images for each
457 of four possible combinations (e.g. Fan In: A1-B1; A2-B1; A4-B4; A5-B4; Fan Out: A3-B3; A3-B3;
458 A6-B5; A6-B6). During reward learning, four B images from each condition were shown (e.g. Fan
459 In: B2 x2; B5 x2; Fan Out: B1 x2; B4 x2) such that two from each condition were paired with
460 reward (e.g. Fan In: B1; Fan Out: B2) and two were paired with a neutral outcome (e.g. Fan In:
461 B4; Fan Out: B5). Finally, during the decision phase, participants made choices between B
462 images that had been directly associated with a reward or neutral outcome (*test* choices) and
463 between A images that had been indirectly associated with these outcomes (*transfer* choices).
464 Test (e.g. Fan In: B1 v B4; Fan Out: B2 v B5) and transfer (e.g. Fan In: A1 v A4; A2 v A5; Fan
465 Out: A3 v A6) choices were made between images from the same condition, and never between
466 images from different conditions.

467 Participants were given a cover story to aid their learning throughout the task. Specifically,
468 participants were told that they were a photographer visiting a new city and would be taking
469 different buses to different locations. At each location, they would be shown a picture they had
470 taken there, and the purpose of the first phase was to learn which photos were taken along each
471 bus route. Then, during the reward learning phase, participants were told that they had returned
472 from their trip and had sent their photos to clients for potential purchase. They were then shown
473 which photos had been purchased and which had not, and their goal was to learn this information.
474 Lastly, during the decision phase, participants were told that they were planning a new trip to the
475 city and were tasked with deciding between bus routes (represented by photos taken on each
476 route) that would take them to locations where they had taken photos their clients purchased.
477 Participants were instructed to use what they had learned (i.e. which photos were taken along the
478 same route and which were or were not purchased) to inform their choices.

479 **MRI Acquisition**

480 MRI data were collected on a 3 T Siemens Magnetom Prisma scanner with a 64-channel head
481 coil. Functional images were acquired using a multiband echo-planar imaging (EPI) sequence
482 (repetition time = 1.5s, echo time = 30ms, flip angle = 65°, acceleration factor = 3, voxel size = 2
483 mm iso, acquisition matrix 96 x 96). Sixty nine oblique axial slices (14° transverse to coronal) were
484 acquired in an interleaved order and spaced 2mm to achieve full brain coverage. Whole-brain
485 high resolution (1 mm iso) T1-weighted structural images were acquired with a magnetization-
486 prepared rapid acquisition gradient-echo (MPRAGE) sequence. Field maps consisting of 69
487 oblique axial slices (2 mm isotropic) were collected to aid registration.

488 **Imaging Data Preprocessing**

489 Results included in this manuscript come from preprocessing performed using *fMRIPrep* 20.2.6,
490 which is based on *Nipype* 1.7.0.⁴⁵

491 **Anatomical Data Preprocessing**

492 Each participant's T1-weighted (T1w) image was corrected for intensity non-uniformity (INU)
493 with N4BiasFieldCorrection⁴⁶, distributed with ANTs 2.3.3⁴⁷ and used as a reference image
494 throughout the workflow. The reference image was then skull-stripped with
495 a *Nipype* implementation of the antsBrainExtraction.sh workflow (from ANTs), using
496 OASIS30ANTs as target template. Brain tissue segmentation of cerebrospinal fluid (CSF), white-
497 matter (WM) and gray-matter (GM) was performed on the brain-extracted T1w using fast⁴⁸ (FSL
498 5.0.9). Volume-based spatial normalization to the *ICBM 152 Nonlinear Asymmetrical template*

499 *version 2009c* (MNI152NLin2009cAsym) standard space was performed through nonlinear
500 registration with *antsRegistration* (ANTs 2.3.3), using brain-extracted versions of both the T1w
501 reference and the T1w template images.

502 **Functional Data Preprocessing**

503 For each of the 5 BOLD runs per subject (two stimulus learning runs, two reward learning runs,
504 and one choice run), the following preprocessing was performed. First, a reference volume and
505 its skull-stripped version were generated using a custom methodology of *fMRIPrep*. A B0-
506 nonuniformity map (or *fieldmap*) was estimated based on two (or more) echo-planar imaging (EPI)
507 references with opposing phase-encoding directions, with *3dQwarp*⁴⁹ (AFNI 20160207). Based
508 on the estimated susceptibility distortion, a corrected EPI reference was calculated for a more
509 accurate co-registration with the anatomical reference. The BOLD reference was then co-
510 registered to the T1w reference using *bbregister* (FreeSurfer) which implements boundary-based
511 registration⁵⁰. Co-registration was configured with six degrees of freedom. Head-motion
512 parameters with respect to the BOLD reference (transformation matrices, and six corresponding
513 rotation and translation parameters) were estimated before any spatiotemporal filtering
514 using *mcflirt*⁵¹ (FSL 5.0.9). BOLD runs were slice-time corrected to 0.708s (0.5 of slice acquisition
515 range 0s-1.42s) using *3dTshift* from AFNI 20160207⁴⁹. The BOLD time-series (including slice-
516 timing correction when applied) were resampled onto their original, native space by applying a
517 single, composite transform to correct for head-motion and susceptibility distortions. The BOLD
518 time-series were resampled into standard space, generating a preprocessed BOLD run in
519 MNI152NLin2009cAsym space. First, a reference volume and its skull-stripped version were
520 generated using a custom methodology of *fMRIPrep*. Several confounding time-series were
521 calculated based on the preprocessed BOLD: framewise displacement (FD), DVARS and three
522 region-wise global signals. FD was computed using two formulations following Power (absolute
523 sum of relative motions)⁵² and Jenkinson (relative root mean square displacement between
524 affines)⁵¹. FD and DVARS are calculated for each functional run, both using their implementations
525 in *Nipype*. The three global signals are extracted within the CSF, the WM, and the whole-brain
526 masks. The head-motion estimates calculated in the correction step were also placed within the
527 corresponding confounds file. The confound time series derived from head motion estimates and
528 global signals were expanded with the inclusion of temporal derivatives and quadratic terms for
529 each⁵³. Frames that exceeded a threshold of 0.5 mm FD or 1.5 standardized DVARS were
530 annotated as motion outliers. All resamplings can be performed with a single interpolation step by
531 composing all the pertinent transformations (i.e. head-motion transform matrices, susceptibility
532 distortion correction when available, and co-registrations to anatomical and output spaces).
533 Gridded (volumetric) resamplings were performed using *antsApplyTransforms* (ANTs),
534 configured with Lanczos interpolation to minimize the smoothing effects of other kernels⁵⁴.
535 Preprocessed data were lastly smoothed using a Gaussian kernel with a FWHM of 6.0mm,
536 masked, and mean-scaled over time.

537 **Functional Imaging Data Analysis**

538 *Beta Series Modeling*

539 Least squares separate (LSS) models were generated for each event (presentation of a category
540 image) in each task following the method described in Turner et al., 2012³² using *Nistats 0.0.1b2*.
541 For each trial, preprocessed data were subjected to a general linear model in which the trial was
542 modeled in its own regressor, while all other trials from that condition were modeled in a second
543 regressor, and other conditions were modeled in their own regressors. Each condition regressor

544 was convolved with the *glover* hemodynamic response function for the model. In addition to
545 condition regressors, 36 nuisance regressors were included in each model consisting of two
546 physiological time series (the mean WM and CSF signals), the global signal, six head-motion
547 parameters, their derivatives, quadratic terms, and squares of derivatives. Spike regression was
548 additionally performed by including a regressor for each motion outlier identified in each run, as
549 in Satterthwaite et al., 2013⁵³. A high-pass filter of 0.0078125 Hz, implemented using a cosine
550 drift model, was also included in each model and AR(1) prewhitening was applied to each model
551 to account for temporal autocorrelation. After fitting each model, the parameter estimate (i.e.,
552 beta) map associated with the target trial's regressor was retained and used for further analysis.
553 Modeling was performed using *NiBetaSeries* 0.6.0⁵⁵ which is based on *Nipype* 1.4.2.⁴⁵ Beta maps
554 for image presentation events, separated by category, for the stimulus learning and reward
555 learning phases and for decisions between images, again separated by category, were used in
556 subsequent analyses.

557 *Multivariate Pattern Decoding Analysis*

558 Beta maps from each trial were next used for multivariate pattern analysis. First, a searchlight
559 classification analysis was conducted for each participant. In brief, a three-way one versus all
560 logistic regression classifier was trained to distinguish categories using leave-one-run-out cross
561 validation from runs of the stimulus learning task. We used winner-take-all labeling to determine
562 the classified label from each trial: the category resulting in the highest probability from the one
563 versus all classification procedure on a given trial was selected as the predicted label for that trial.
564 Input data were selected using a spherical searchlight (radius = 2 voxels) moved around the whole
565 brain. Although the experimental design leads the class labels for each category to be imbalanced
566 during the stimulus learning phase (i.e. one label always has twice as many occurrences as the
567 other two), we dealt with this label imbalance in two ways. First, the class weights applied to each
568 category by the classifier were determined using the 'balanced' keyword in *sklearn*⁵⁶ such that the
569 weights were the number of samples divided by the number of labels (3) multiplied by the total
570 number of occurrences of each label. Second, our metric of performance was the weighted-F1
571 score, which is the harmonic mean of precision and recall. Each of these methods are commonly
572 used in the machine learning literature to deal with class imbalance in training data. For each
573 searchlight sphere, we additionally computed chance performance via a permutation test: labels
574 were shuffled 1000 times and the weighted F1-score resulting from each of these permutations
575 was computed. Chance classification performance was then calculated as the 95th percentile of
576 the F1-score permutation distribution. For each voxel, we then subtracted chance level
577 performance from the classification accuracy to produce a map of corrected classification
578 performance for each participant. Finally, an FDR-corrected ($q < 0.05$) group-level map over all
579 individual subject difference maps was created.

580 Following classifier training on the stimulus learning phase, we then tested the classifier on runs
581 from both the reward learning and decision phases. Functional data from each participant on each
582 of these phases of the experiment was first masked using the group-level searchlight map
583 produced from the previously described procedure. The three-way logistic regression classifier
584 was then re-trained on both runs of the stimulus learning phase, using only these voxels, and then
585 tested separately on the reward learning and decision phases. L1-regularization was used to
586 reduce overfitting in this procedure. We again used the weighted F1-score as our accuracy metric,
587 and the 95th percentile of the permutation distribution as our measure of chance classifier
588 performance.

589 Finally, to address our primary question, we created an index of memory reactivation from the
590 classifier. Specifically, for each trial, we extracted the probability that the classifier assigned to
591 each category label. A trial was then considered a trial on which memory reactivation occurred if
592 the following criteria were met: i) the true category label was assigned the highest probability by
593 the classifier and ii) the associated category was assigned the second highest probability by the
594 classifier. If these criteria were met, the trial was assigned a one and, if not, a zero. Our logic for
595 using this criteria was conservative: we reasoned that the classifier should always assign the
596 highest probability to the category represented by the image that is presently shown on the
597 screen. Because, by definition, both off-screen categories were candidates for association when
598 presented as part of Fan In trials during the reward learning phase, we were unable to calculate
599 a reactivation score for these trials. We were further limited in our ability to compare reactivation
600 across phases because the classifier was more accurate at identifying category images presented
601 during the decision phase than during the reward learning phase. We were, however, able to
602 investigate individual differences in reactivation for Fan Out trials between phases by accounting
603 for this difference in classification performance by z-scoring reactivation scores within each
604 phase, as this removes group-level differences while leaving individual differences intact.

605 **Regression Analyses**

606 Unless otherwise noted, parameters for all regression models described here were estimated
607 using hierarchical Bayesian inference such that group-level priors were used to regularize subject-
608 level estimates. The joint posterior was approximated using No-U-Turn Sampling⁵⁷ as
609 implemented in stan. Four chains with 2000 samples (1000 discarded as burn-in) were run for a
610 total of 4000 posterior samples per model. Chain convergence was determined by ensuring that
611 the Gelman-Rubin statistic \hat{R} was close to 1. Default weakly-informative priors implemented in the
612 *rstanarm*⁵⁸ package were used for each regression model. For all models, fixed effects are
613 reported in the text as the mean of each parameter's marginal posterior distribution alongside
614 95% or 90% credible intervals, which indicate where that percentage of the posterior density falls.
615 Parameter values outside of this range are unlikely given the model, data, and priors. Thus, if the
616 range of likely values does not include zero, we conclude that a meaningful effect was observed.

617 We first assessed choice performance on the decision phase of the task. For each subject s and
618 trial t , a mixed effects logistic regression was used to predict if the correct image was chosen:

$$619 \quad p(\text{Correct}_t) = \sigma(\beta_0 + b_{0,s[t]} + \text{Condition}_t(\beta_1 + b_{1,s[t]}))$$

$$620 \quad \sigma(x) = \frac{1}{1 + e^{-x}}$$

621 where *Correct* was equal to 1 if the participant chose either the image directly associated with
622 reward (in the case of test trials) or the image indirectly associated with reward (in the case of
623 transfer trials), and *Condition* was a categorical variable coded as 0.5 for Fan In trials and -0.5
624 for Fan Out trials. This model was fit separately for test and transfer choices.

625 We also assessed the relationship between response time and accuracy during transfer choices
626 using the following mixed effects logistic regression, which included an additional main effect of
627 response time as well the interaction between response time and condition:

$$628 \quad p(\text{Correct}_t) = \sigma(\beta_0 + b_{0,s[t]} + \text{Condition}_t * (\beta_1 + b_{1,s[t]}) + RT_t * (\beta_2 + b_{2,s[t]}) \\ 629 \quad + \text{Condition}_t \times RT_t * (\beta_3 + b_{3,s[t]}))$$

630 where RT was the response time on each transfer choice trial.

631 We determined the ability of the trained MVPA classifier to distinguish each category label from
632 chance using the following mixed effects linear regression:

$$633 \quad Accuracy - Chance = \beta_0 + b_{0,s[t]} + Phase_t(\beta_1 + b_{1,s[t]})$$

634 where $Accuracy - Chance$ was the 95th percentile of the permutation distribution subtracted from
635 classification accuracy, and $Phase$ was a categorical variable coded as 0.5 for the decision phase
636 and -0.5 for the reward learning phase. This model was fit separately for each category (face,
637 scene and object).

638 Another set of models was fit to assess the relationship between memory reactivation and transfer
639 choice behavior. Analyses were conducted on the average reactivation level of each stimulus. In
640 order to assess effects of reactivation on transfer accuracy for each stimulus, i , accuracy was first
641 transformed⁵⁹ to ensure that all responses fell within the interval (0,1):

$$642 \quad TransAcc'_i = \frac{TransAcc_i(N - 1) + 0.5}{N}$$

643 where $TransAcc$ was participants' average transfer accuracy for each consequent stimulus and
644 N was the sample size (39). We first examined the effect of (z-scored) differences in reactivation
645 between the reward learning and decision phases for each associated antecedent-consequent
646 pair of Fan Out stimuli on transfer accuracy. To do so, we fit a mixed effects beta regression:

$$647 \quad \text{logit}(TransAcc'_i) = \beta_0 + b_{0,s[i]} + \Delta Reactivation_t(\beta_1 + b_{1,s[i]})$$

648 where $\Delta Reactivation$ is the difference in memory reactivation between reward learning and the
649 decision phase for each pair. Similar beta regressions were used to assess effects of memory
650 reactivation during the decision phase for Fan In and Fan Out consequent stimuli, separately. To
651 assess effects on choice transfer response time, linear mixed effects regressions with the same
652 predictors were used instead.

653 We additionally assessed how memory reactivation differed for each condition (Fan In or Fan Out)
654 during the decision phase. We performed this analysis using the following mixed effects linear
655 regression:

$$656 \quad Reactivation = \beta_0 + b_{0,s} + Condition (\beta_1 + b_{1,s})$$

657 where $Reactivation$ was memory reactivation during the decision phase for each participant and
658 condition and $Condition$ was coded identically to the models described above.

659 Lastly, we examined individual differences in strategy usage by comparing our reactivation
660 measures across phases of the task. Specifically, we fit a simple linear regression predicting each
661 participants' average level of memory reactivation for Fan Out during reward learning from their
662 difference in memory reactivation during the decision phase.

663 **Acknowledgements**

664 The authors thank members of the Shohamy Lab for insightful discussion. Support was provided
665 by the NSF GRFP (J.N.; award #1644869), the NSF (D.S., N.D.; award #1822619), the NIMH/NIH

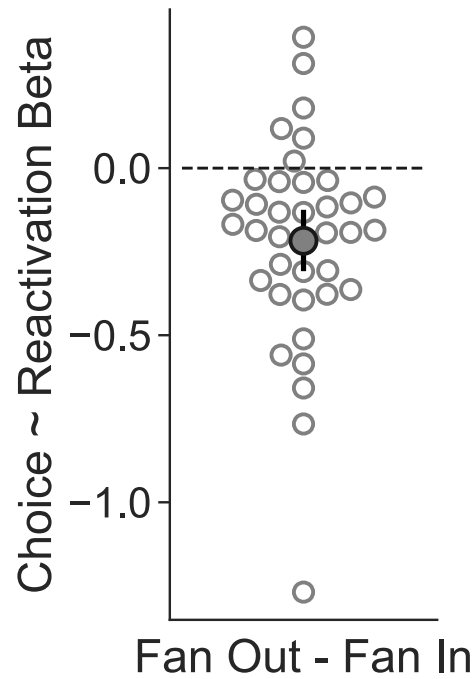
666 (D.S., N.D.; award #MH121093) and the Templeton Foundation (D.S. grant # 60844). Figures 1
667 and 3 use images created by Creative, Freepik, and Paul J on Flaticon.com.

668 **References**

- 669 1. Sutton, R. S. Integrated Architectures for Learning, Planning, and Reacting Based on
670 Approximating Dynamic Programming. in *Machine Learning Proceedings 1990* (eds. Porter,
671 B. & Mooney, R.) 216–224 (Morgan Kaufmann, 1990). doi:10.1016/B978-1-55860-141-
672 3.50030-4.
- 673 2. Mattar, M. G. & Lengyel, M. Planning in the brain. *Neuron* **110**, 914–934 (2022).
- 674 3. LaValle, S. M. *Planning Algorithms*. (Cambridge University Press, 2006).
675 doi:10.1017/CBO9780511546877.
- 676 4. Daw, N. D. & Dayan, P. The algorithmic anatomy of model-based evaluation. *Philos. Trans.*
677 *R. Soc. B Biol. Sci.* **369**, 20130478–20130478 (2014).
- 678 5. Doll, B. B., Simon, D. A. & Daw, N. D. The ubiquity of model-based reinforcement learning.
679 *Curr. Opin. Neurobiol.* **22**, 1075–1081 (2012).
- 680 6. Dolan, R. J. & Dayan, P. Goals and Habits in the Brain. *Neuron* **80**, 312–325 (2013).
- 681 7. O’Keefe, J. & Nadel, L. *The hippocampus as a cognitive map*. (Clarendon Press ; Oxford
682 University Press, 1978).
- 683 8. Tolman, E. C. Cognitive maps in rats and men. *Psychol. Rev.* **55**, 189–208 (1948).
- 684 9. Daw, N. D., Gershman, S. J., Seymour, B., Dayan, P. & Dolan, R. J. Model-Based Influences
685 on Humans’ Choices and Striatal Prediction Errors. *Neuron* **69**, 1204–1215 (2011).
- 686 10. Doll, B. B., Duncan, K. D., Simon, D. A., Shohamy, D. & Daw, N. D. Model-based choices
687 involve prospective neural activity. *Nat. Neurosci.* **18**, 767–772 (2015).
- 688 11. Wimmer, G. E. & Shohamy, D. Preference by Association: How Memory Mechanisms in the
689 Hippocampus Bias Decisions. *Science* **338**, 270–273 (2012).
- 690 12. Kurth-Nelson, Z., Barnes, G., Sejdinovic, D., Dolan, R. & Dayan, P. Temporal structure in
691 associative retrieval. *eLife* **4**, e04919 (2015).
- 692 13. Shohamy, D. & Wagner, A. D. Integrating Memories in the Human Brain: Hippocampal-
693 Midbrain Encoding of Overlapping Events. *Neuron* **60**, 378–389 (2008).
- 694 14. Jones, J. L. *et al.* Orbitofrontal Cortex Supports Behavior and Learning Using Inferred But Not
695 Cached Values. *Science* **338**, 953–956 (2012).
- 696 15. Mattar, M. G. & Daw, N. D. Prioritized memory access explains planning and hippocampal
697 replay. *Nat. Neurosci.* **21**, 1609–1617 (2018).
- 698 16. Brogden, W. J. Sensory pre-conditioning. *J. Exp. Psychol.* **25**, 323 (19400101).
- 699 17. Diba, K. & Buzsáki, G. Forward and reverse hippocampal place-cell sequences during ripples.
700 *Nat. Neurosci.* **10**, 1241–1242 (2007).
- 701 18. Johnson, A. & Redish, A. D. Neural Ensembles in CA3 Transiently Encode Paths Forward of
702 the Animal at a Decision Point. *J. Neurosci.* **27**, 12176–12189 (2007).
- 703 19. Pfeiffer, B. E. & Foster, D. J. Hippocampal place-cell sequences depict future paths to
704 remembered goals. *Nature* **497**, 74–79 (2013).
- 705 20. Gruber, M. J., Ritchey, M., Wang, S.-F., Doss, M. K. & Ranganath, C. Post-learning
706 Hippocampal Dynamics Promote Preferential Retention of Rewarding Events. *Neuron* **89**,
707 1110–1120 (2016).
- 708 21. Murty, V. P., Tompary, A., Adcock, R. A. & Davachi, L. Selectivity in Postencoding
709 Connectivity with High-Level Visual Cortex Is Associated with Reward-Motivated Memory. *J.*
710 *Neurosci.* **37**, 537–545 (2017).
- 711 22. Gershman, S. J., Markman, A. B. & Otto, A. R. Retrospective revaluation in sequential
712 decision making: A tale of two systems. *J. Exp. Psychol. Gen.* **143**, 182–194 (2014).

- 713 23. Momennejad, I., Otto, A. R., Daw, N. D. & Norman, K. A. Offline replay supports planning in
714 human reinforcement learning. *eLife* **7**, e32548 (2018).
- 715 24. Foster, D. J. & Wilson, M. A. Reverse replay of behavioural sequences in hippocampal place
716 cells during the awake state. *Nature* **440**, 680–683 (2006).
- 717 25. Michon, F., Sun, J.-J., Kim, C. Y., Ciliberti, D. & Kloosterman, F. Post-learning Hippocampal
718 Replay Selectively Reinforces Spatial Memory for Highly Rewarded Locations. *Curr. Biol.* **29**,
719 1436–1444.e5 (2019).
- 720 26. Singer, A. C. & Frank, L. M. Rewarded Outcomes Enhance Reactivation of Experience in the
721 Hippocampus. *Neuron* **64**, 910–921 (2009).
- 722 27. Liu, Y., Mattar, M. G., Behrens, T. E. J., Daw, N. D. & Dolan, R. J. Experience replay is
723 associated with efficient nonlocal learning. *Science* **372**, eabf1357 (2021).
- 724 28. Shohamy, D. & Daw, N. D. Integrating memories to guide decisions. *Curr. Opin. Behav. Sci.*
725 **5**, 85–90 (2015).
- 726 29. Cohen, N. J. & Eichenbaum, H. *Memory, amnesia, and the hippocampal system*. xii, 330 (The
727 MIT Press, 1993).
- 728 30. Haxby, J. V. *et al.* Distributed and Overlapping Representations of Faces and Objects in
729 Ventral Temporal Cortex. *Science* **293**, 2425–2430 (2001).
- 730 31. Polyn, S. M., Natu, V. S., Cohen, J. D. & Norman, K. A. Category-Specific Cortical Activity
731 Precedes Retrieval During Memory Search. *Science* **310**, 1963–1966 (2005).
- 732 32. Turner, B. O., Mumford, J. A., Poldrack, R. A. & Ashby, F. G. Spatiotemporal activity
733 estimation for multivoxel pattern analysis with rapid event-related designs. *NeuroImage* **62**,
734 1429–1438 (2012).
- 735 33. Gupta, A. S., van der Meer, M. A. A., Touretzky, D. S. & Redish, A. D. Hippocampal Replay
736 Is Not a Simple Function of Experience. *Neuron* **65**, 695–705 (2010).
- 737 34. Ambrose, R. E., Pfeiffer, B. E. & Foster, D. J. Reverse Replay of Hippocampal Place Cells Is
738 Uniquely Modulated by Changing Reward. *Neuron* **91**, 1124–1136 (2016).
- 739 35. Moore, A. W. & Atkeson, C. G. Prioritized sweeping: Reinforcement learning with less data
740 and less time. *Mach. Learn.* **13**, 103–130 (1993).
- 741 36. Coulom, R. Efficient Selectivity and Backup Operators in Monte-Carlo Tree Search. in
742 *Computers and Games* (eds. van den Herik, H. J., Ciancarini, P. & Donkers, H. H. L. M.) vol.
743 4630 72–83 (Springer Berlin Heidelberg, 2007).
- 744 37. Staresina, B. P., Alink, A., Kriegeskorte, N. & Henson, R. N. Awake reactivation predicts
745 memory in humans. *Proc. Natl. Acad. Sci.* **110**, 21159–21164 (2013).
- 746 38. Tambini, A. & Davachi, L. Persistence of hippocampal multivoxel patterns into postencoding
747 rest is related to memory. *Proc. Natl. Acad. Sci.* **110**, 19591–19596 (2013).
- 748 39. Lee, A. K. & Wilson, M. A. Memory of Sequential Experience in the Hippocampus during Slow
749 Wave Sleep. *Neuron* **36**, 1183–1194 (2002).
- 750 40. Dupret, D., O'Neill, J., Pleydell-Bouverie, B. & Csicsvari, J. The reorganization and
751 reactivation of hippocampal maps predict spatial memory performance. *Nat. Neurosci.* **13**,
752 995–1002 (2010).
- 753 41. Ólafsdóttir, H. F., Barry, C., Saleem, A. B., Hassabis, D. & Spiers, H. J. Hippocampal place
754 cells construct reward related sequences through unexplored space. *eLife* **4**, e06063 (2015).
- 755 42. Sharp, M. E., Foerde, K., Daw, N. D. & Shohamy, D. Dopamine selectively remediates 'model-
756 based' reward learning: a computational approach. *Brain* **139**, 355–364 (2016).
- 757 43. Sharpe, M. J. *et al.* Dopamine transients are sufficient and necessary for acquisition of model-
758 based associations. *Nat. Neurosci.* **20**, 735–742 (2017).
- 759 44. Sharp, P. B. & Eldar, E. Humans Adaptively Deploy Forward and Backward Prediction. (2023)
760 doi:10.31234/osf.io/wdbg4.

- 761 45. Gorgolewski, K. *et al.* Nipype: A Flexible, Lightweight and Extensible Neuroimaging Data
762 Processing Framework in Python. *Front. Neuroinformatics* **5**, (2011).
- 763 46. Tustison, N. J. *et al.* N4ITK: Improved N3 Bias Correction. *IEEE Trans. Med. Imaging* **29**,
764 1310–1320 (2010).
- 765 47. Avants, B. B., Epstein, C. L., Grossman, M. & Gee, J. C. Symmetric diffeomorphic image
766 registration with cross-correlation: Evaluating automated labeling of elderly and
767 neurodegenerative brain. *Med. Image Anal.* **12**, 26–41 (2008).
- 768 48. Zhang, Y., Brady, M. & Smith, S. Segmentation of brain MR images through a hidden Markov
769 random field model and the expectation-maximization algorithm. *IEEE Trans. Med. Imaging*
770 **20**, 45–57 (2001).
- 771 49. Cox, R. W. & Hyde, J. S. Software tools for analysis and visualization of fMRI data. *NMR*
772 *Biomed.* **10**, 171–178 (1997).
- 773 50. Greve, D. N. & Fischl, B. Accurate and robust brain image alignment using boundary-based
774 registration. *NeuroImage* **48**, 63–72 (2009).
- 775 51. Jenkinson, M., Bannister, P., Brady, M. & Smith, S. Improved Optimization for the Robust and
776 Accurate Linear Registration and Motion Correction of Brain Images. *NeuroImage* **17**, 825–
777 841 (2002).
- 778 52. Power, J. D. *et al.* Methods to detect, characterize, and remove motion artifact in resting state
779 fMRI. *NeuroImage* **84**, 320–341 (2014).
- 780 53. Satterthwaite, T. D. *et al.* An improved framework for confound regression and filtering for
781 control of motion artifact in the preprocessing of resting-state functional connectivity data.
782 *NeuroImage* **64**, 240–256 (2013).
- 783 54. Lanczos, C. Evaluation of Noisy Data. *J. Soc. Ind. Appl. Math. Ser. B Numer. Anal.* **1**, 76–85
784 (1964).
- 785 55. Kent, J. & Herholz, P. NiBetaSeries: task related correlations in fMRI. *J. Open Source Softw.*
786 **4**, 1295 (2019).
- 787 56. Pedregosa, F. *et al.* Scikit-learn: Machine Learning in Python. *J. Mach. Learn. Res.* **12**, 2825–
788 2830 (2011).
- 789 57. Hoffman, M. D. & Gelman, A. The No-U-Turn Sampler: Adaptively Setting Path Lengths in
790 Hamiltonian Monte Carlo. 31.
- 791 58. Bayesian Applied Regression Modeling via Stan. <https://mc-stan.org/rstanarm/>.
- 792 59. Smithson, M. & Verkuilen, J. A better lemon squeezer? Maximum-likelihood regression with
793 beta-distributed dependent variables. *Psychol. Methods* **11**, 54 (20060403).
- 794



795

796 **Supplementary Figure 1.** As shown in Figure 4C, greater memory reactivation at decision time is
797 associated with less effective transfer decisions for Fan In but not Fan Out image pairs. Shown here is this
798 effect, the difference in slopes, for every participant and at the group-level. Participants demonstrate this
799 relationship more for Fan Out than Fan In decisions ($\beta_0 = -0.215$, 95% $CI = [-0.312, -0.118]$), as
800 indicated by comparing their random slopes. The filled point represents the group-average difference in
801 slopes, whereas empty points represent individual slope differences.