

1 When to retrieve and encode 2 episodic memories: a neural 3 network model of 4 hippocampal-cortical interaction

5 Qihong Lu^{1,2*}, Uri Hasson^{1,2}, Kenneth A. Norman^{1,2}

*For correspondence:
qlu@princeton.edu

6 ¹Department of Psychology, Princeton University; ²Princeton Neuroscience Institute,
7 Princeton University

8

9 **Abstract** Recent human behavioral and neuroimaging results suggest that people are selective
10 in when they encode and retrieve episodic memories. To explain these findings, we trained a
11 memory-augmented neural network to use its episodic memory to support prediction of
12 upcoming states in an environment where past situations sometimes reoccur. We found that the
13 network learned to retrieve selectively as a function of several factors, including its uncertainty
14 about the upcoming state. Additionally, we found that selectively encoding episodic memories at
15 the end of an event (but not mid-event) led to better subsequent prediction performance. In all of
16 these cases, the benefits of selective retrieval and encoding can be explained in terms of
17 reducing the risk of retrieving irrelevant memories. Overall, these modeling results provide a
18 resource-rational account of why episodic retrieval and encoding should be selective and lead to
19 several testable predictions.

20

21 Introduction

22 In a natural setting, when should an intelligent agent encode and retrieve episodic memories? For
23 example, suppose I am viewing the BBC television series *Sherlock*. Should I retrieve an episodic
24 memory that I formed when I watched earlier parts of the show, and if so, when should I retrieve
25 this memory? When should I encode information about the ongoing episode?

26 Although episodic memory is one of the most studied topics in cognitive psychology and cogni-
27 tive neuroscience, the answers to these questions are still unclear, in large part because episodic
28 memory research has traditionally focused on experiments using simple, well-controlled stimuli,
29 where participants receive clear instructions about when to encode and retrieve. For example, a
30 typical episodic memory experiment could ask participants to remember a set of random word-
31 pairs; later on, given a word-cue, the participants need to report the associated word (*Kahana,*
32 *2012*). In this kind of word-pair experiment, the optimal timing for encoding and retrieval is clear:
33 The participant should encode an episodic memory when they study a word-pair and retrieve the
34 associate when they are prompted by a cue. Existing computational models of human memory
35 have similarly focused on discretized list-learning paradigms like the (hypothetical) word-pair learn-
36 ing study described above – these models (see *Norman et al. 2008* for a review) are primarily de-
37 signed to answer questions about what happens as a result of a particular sequence of encoding
38 and retrieval trials, not questions about when encoding and retrieval should occur in the first place.

39 Recently, there has been increasing interest in using naturalistic stimuli such as movies or au-

40 dio narratives in psychological experiments, to complement results from traditional experiments
41 using simple and well-controlled stimuli (*Sonkusare et al., 2019; Nastase et al., 2020*). These ex-
42 periments have the potential to shed light on when encoding and retrieval take place during event
43 perception in a naturalistic context, where no one is explicitly instructing participants about how
44 to use episodic memory. These studies have found evidence that episodic encoding and retrieval
45 occur *selectively* over time. For example, results from fMRI studies suggest that episodic encoding
46 occurs preferentially at the ends of events (*Baldassano et al., 2017; Ben-Yakov et al., 2013; Ben-*
47 *Yakov and Henson, 2018; Reagh et al., 2020*), and episodic retrieval happens preferentially when
48 people are uncertain about the ongoing situation (*Chen et al., 2016*). Selectivity effects can also
49 be observed in the realm of more traditional list-learning studies – for example, there is exten-
50 sive behavioral and neuroscientific evidence that stimuli that trigger strong prediction errors are
51 preferentially encoded into episodic memory (for reviews, see *Frank and Kafkas 2021; Quent et al.*
52 *2021b*).

53 The goal of the present work is to develop a computational model that can account for *when*
54 *episodic encoding and retrieval take place* in naturalistic situations; the model is meant to capture
55 key features of cortical-hippocampal interactions, as described below. We formalize the task of
56 event processing by assuming that events involve sequences of states drawn from some underly-
57 ing event schema, and that the agent’s goal is to predict upcoming states. We then seek to identify
58 policies for episodic encoding and retrieval by optimizing a neural network model on the event pro-
59 cessing task. We analyze how the optimal policy changes under different environmental regimes,
60 and how well this policy captures human behavioral and neuroimaging data. To the extent that
61 they match, the model can be viewed as providing a *resource-rational* account of those findings
62 (i.e., an explanation of how these encoding and retrieval policies arise as a joint adaptation to the
63 constraints imposed by the human cognitive architecture and the constraints imposed by the task
64 environment; *Griffiths et al. 2015; Lieder and Griffiths 2019*; see also *Anderson and Schooler 2000;*
65 *Gershman 2021*).

66 Overall, we find that the best-performing policies are selective in when encoding and retrieval
67 take place, and that the types of selectivity identified by the model line up well with types of selec-
68 tivity identified empirically. The key intuition behind these effects is that – while retrieving episodic
69 memories can help us to predict upcoming states – there are risks to episodic retrieval: If you
70 retrieve an irrelevant memory, you could make confident, wrong predictions that have negative
71 consequences. The selective encoding and retrieval policies identified by the model help it to mit-
72 igate these risks while retaining the benefits of episodic memory. In the sections that follow, we
73 describe our cortical-hippocampal model, how we applied it to the tasks of interest, and the results
74 of our simulations.

75 **A neural network model of cortical-hippocampal interaction**

76 Our modeling work leverages recent advances in memory-augmented neural networks (*Graves*
77 *et al., 2016; Ritter et al., 2018*), deep reinforcement learning (*Mnih et al., 2016; Sutton and Barto,*
78 *2018*), and meta-learning (*Wang et al., 2018; Botvinick et al., 2019*) – these advances (collectively)
79 make it possible for neural network models to *learn to use episodic memory* in the service of predic-
80 tion.

81 Our model (Figure 1A) has two parts, which are meant to correspond to cortex and hippocam-
82 pus, and which collectively implement three key memory systems (working memory, semantic
83 memory, and episodic memory). The cortical part of the model incorporates a Long-Short-Term
84 Memory module (LSTM; *Hochreiter and Schmidhuber 1997*), which is a recurrent neural network
85 (RNN) with gating mechanisms. In addition to the LSTM module, the cortical network also incor-
86 porates a nonlinear decision layer (to assist with mapping inputs to next-state predictions) and an
87 episodic memory (EM) gating layer, the function of which is described below. The LSTM module
88 gives the cortical network the ability to actively maintain and integrate information over time. For
89 terminological convenience, we will refer to this active maintenance ability in the paper as “working

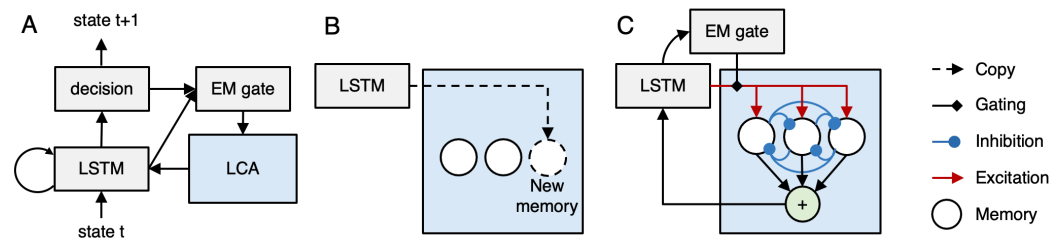


Figure 1. Cortical-hippocampal Model. A) At a given moment, the cortical part of the model (shown in gray) observes the current state and predicts the upcoming state. It incorporates a Long Short Term Memory (LSTM; *Hochreiter and Schmidhuber, 1997*) network, which integrates information over time; the LSTM feeds into a non-linear decision layer. The LSTM and decision layers also project to an episodic memory (EM) gating layer that determines when episodic memories are retrieved (see part C of figure). The entire cortical network is trained by an advantage actor critic (A2C) objective (*Mnih et al., 2016*) to optimize next-state prediction. B) Episodic memory encoding involves copying the current hidden state and appending it to the list of memories stored in the episodic memory system (shown in blue), which is meant to correspond to hippocampus. C) Episodic memory retrieval is implemented using a leaky competing accumulator model (LCA; *Usher and McClelland, 2001*) – each memory receives excitation proportional to its similarity to the current hidden state, and different memories compete with each other via lateral inhibition. The EM gate (whose value is set by the EM gate layer of the cortical network) scales the level of excitation coming into the network. After a fixed number of time steps, an activation-weighted sum of all memories is added back to the cell state of the LSTM.

90 memory". However, we should emphasize that – contrary to classic views of working memory (e.g.,
 91 *Baddeley 2000*) – our model does not have a working memory buffer that is set apart from other
 92 parts of the model that do stimulus processing; rather, active maintenance and integration are
 93 accomplished via recurrent activity in the parts of the model that are doing stimulus processing.
 94 In this respect, the architecture of our model fits with the *process memory* framework set forth by
 95 *Hasson et al. (2015)*. In addition to this active maintenance ability, the connection weights of the
 96 cortical network gradually extract regularities from the environment over time; this gradual learn-
 97 ing of regularities can be viewed as an implementation of semantic memory (*Rumelhart et al.,*
 98 *1987; McClelland and Rogers, 2003; Rogers and McClelland, 2004; Saxe et al., 2019*).

99 The cortical network is also connected to an episodic memory module (meant to simulate hip-
 100 pocampus) that stores snapshots of cortical activity patterns (Figure 1B) and reinstates these pat-
 101 terns to the cortical network; see the next section for more information on the model's encoding
 102 policy (i.e., when it stores snapshots). Episodic memory retrieval (Figure 1C) is implemented via a
 103 leaky competing accumulator process (LCA; *Usher and McClelland 2001; Polyn et al. 2009*). In the
 104 LCA, memories compete to be retrieved according to how well they match the current state of the
 105 cortical network, and the output of this competitive retrieval process is added back into the corti-
 106 cal network. Crucially, the degree to which memories are activated during the retrieval process is
 107 multiplicatively gated by the EM gate layer of the cortical network – this gives the cortical network
 108 the ability to shape when episodic retrieval occurs (for more details on how EM works in the model,
 109 see the *Episodic retrieval* section in the *Methods*).

110 The entire cortical network (composed of the LSTM, decision, and EM gate layers) is trained
 111 via a reinforcement learning algorithm to optimize prediction of the next state given the current
 112 state as input; the trainable nature of the EM gate allows the network to learn a policy for when
 113 episodic memory retrieval should occur, in order to optimize next-state prediction. Specifically, we
 114 used a meta-learning procedure (*Wang et al., 2018*) whereby the model was trained repeatedly on
 115 all conditions of interest with modifiable cortical weights (*meta-training*), before being evaluated in
 116 these conditions with cortical weights frozen (*meta-testing*). This procedure captures the idea that
 117 cortical weights only change gradually (*McClelland et al., 1995*), and thus are unlikely to be modified
 118 enough by one experience to support recall of unique aspects of that experience; as such, memory
 119 for these unique details depends critically on that information being held in working memory or
 120 episodic memory (for more details, see the *Model training and testing* section in the *Methods*).

121 During meta-training, the model is rewarded for correct next-state predictions and punished for
122 incorrect next-state predictions; we also gave the model the option of saying “don’t know” (instead
123 of predicting a specific next state), in which case it receives zero reward. In the real world, there are
124 often different costs associated with making commission errors (wrong predictions) and omission
125 errors (not making a prediction). Having the “don’t know” option gives the model the freedom
126 to choose whether it should make a specific prediction (thereby incurring the risk of making a
127 commission error and receiving a penalty) or whether it should express uncertainty to avoid a
128 possible penalty. Intuitively, this choice should depend on the environment. For example, if the
129 penalty for misprediction is zero, the model should make a prediction even if it has high uncertainty
130 about the upcoming state. In contrast, if the penalty for misprediction is high, the model should
131 only make a prediction if it is certain about what would happen next. Practically speaking, the
132 consequence of including the “don’t know” option is to induce the model to wait longer to retrieve
133 episodic memories (see results below and also Appendix 5).

134 **Modeling the contribution of episodic memory to naturalistic event under-** 135 **standing**

136 Our initial modeling target was a recent study by *Chen et al. (2016)*, which explored the role of
137 episodic memory in naturalistic event understanding. In this study, participants viewed an episode
138 from the *Twilight Zone* television series. This episode was divided into two parts (part 1 and part 2).
139 Participants in the recent memory (RM) condition viewed the two parts back-to-back; participants in
140 the distant memory (DM) condition had a one-day gap in between the two parts of this TV episode;
141 participants in the no memory (NM) condition only watched the second part (*Chen et al., 2016*). In
142 the RM condition, participants can build up a situation model – i.e., a representation of the relevant
143 features of the ongoing situation (*Richmond and Zacks, 2017; Stawarczyk et al., 2019; Zacks, 2020;*
144 *Ranganath and Ritchey, 2012*) – during the first part of the movie and actively maintain it over
145 time; all of that information is still actively represented at the start of part 2. By contrast, in the
146 DM condition, a day has passed between part 1 and part 2, so participants are no longer actively
147 maintaining the relevant situation model at the start of part 2.

148 Taken together, these conditions can be viewed as manipulating the *availability* of relevant
149 episodic memories and also the *demand* for episodic retrieval. In the NM condition, at the start
150 of part 2, participants have gaps in their situation model (because they did not view part 1) and
151 thus there is a strong demand to fill those gaps, to better understand what is going on; however,
152 they do not have any relevant episodic memories available to fill those gaps. In the DM condition,
153 because of the one-day delay, participants also have gaps in their representation of the situation
154 in working memory that need to be filled with information from part 1; however, unlike the NM
155 participants, DM participants can meet this demand by retrieving information about part 1 that
156 was stored in episodic memory. In the RM condition, like the DM condition, participants have rel-
157 evant information about part 1 available in episodic memory (participants’ experience in part 1 of
158 the DM and RM conditions was identical, so presumably they stored the same episodic memories
159 during part 1), but there is less of a demand to retrieve these episodic memories in the RM con-
160 dition (because these participants were not interrupted, and thus these participants should have
161 fewer gaps in their understanding of the situation). The comparison of the RM and DM conditions
162 is thus a relatively pure manipulation of demand for episodic memory retrieval. If episodic mem-
163 ory retrieval is sensitive to the need to retrieve (i.e., whether there are gaps to fill in), then more
164 retrieval should take place in the DM condition, but if episodic memory retrieval is automatic, re-
165 trieval should occur at similar levels in the RM and DM conditions. The results of the *Chen et al.*
166 *(2016)* study strongly support the former (“demand-sensitive”) view of episodic retrieval. During
167 the first two minutes of part 2, the researchers found strong hippocampal-cortical activity coupling
168 measured using inter-subject functional connectivity (ISFC; *Simony et al. 2016*) for DM participants,
169 while the level of coupling was much weaker for participants in the RM and NM conditions (*Chen*

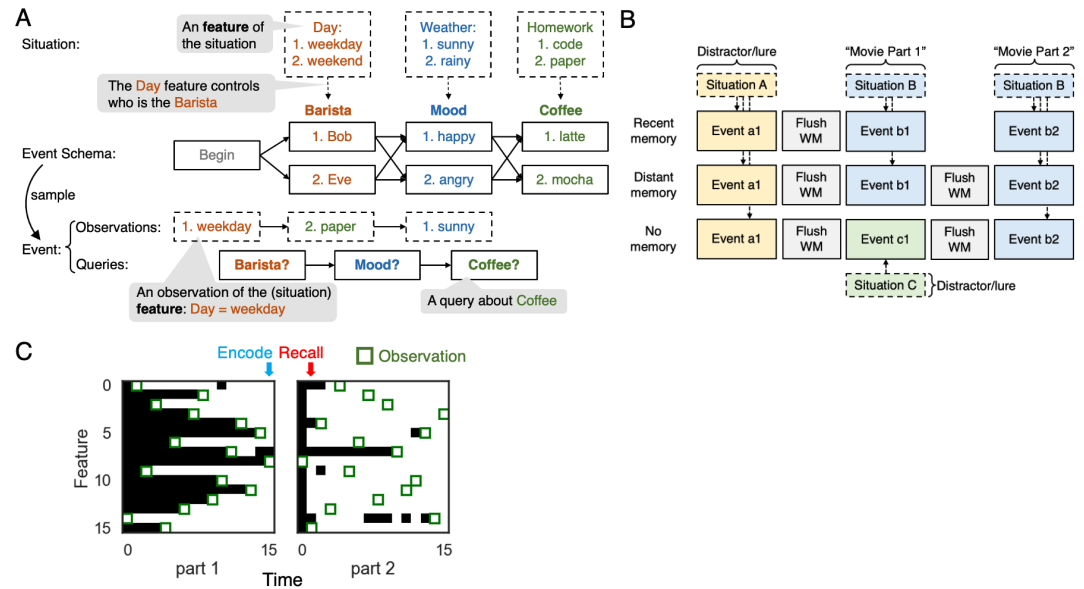


Figure 2. A situation-dependent event processing task. A) An *event* is a sequence of states, sampled from an event schema and conditioned on a situation. An *event schema* is a graph where each node is a state. A *situation* is a collection of features (e.g., Day, Weather, Homework) set to particular values (e.g., Day = weekday). The features of the current situation control how the event unfolds (e.g., the value of the Day feature controls which Barista state is observed). At each time point, the network observes the value of a randomly selected feature of the current situation, and responds to a query about what will happen next. B) We created three task conditions to simulate the design used by *Chen et al. (2016)*: recent memory (RM), distant memory (DM), and no memory (NM); see text for details. C) Decoded contents of the model's working memory for an example trial from the DM condition. Green boxes indicate time points where the value of a particular situation feature was observed. The color of a square indicates whether the correct (i.e., observed) value of that feature can be decoded from the model's working memory state (white = accurate decoding; black = inaccurate decoding). See text for additional explanation.

170 *et al., 2016*). Notably, cortical regions that had a strong coupling with the hippocampus (in the
 171 DM condition) largely overlapped with the default mode network (DMN), which is believed to ac-
 172 tively maintain a situation model (*Stawarczyk et al., 2019*). These results fit with the idea that more
 173 information is being communicated between hippocampus and cortex in the DM (“high episodic
 174 memory demand”) condition than in the RM (“low episodic memory demand”) condition and the
 175 NM condition (where there are no relevant episodic memories to retrieve). This “demand sensitive”
 176 view of episodic memory implies that cortex can be strategic in how it calls upon the hippocampus
 177 to support event understanding, and it underlines the importance of the aforementioned goal of
 178 characterizing the policy for when retrieval should occur.

179 Training environment

180 To simulate the task of event processing, we define an *event* as a sequence of states, sampled from
 181 an underlying graph that represents the *event schema*. Figure 2A shows a “coffee shop visit” event
 182 schema graph with three time points; each time point has two possible states. Each instance of an
 183 event (here, each visit to the coffee shop) is associated with a *situation* – a collection of features set
 184 to particular values; importantly, the features of the current situation determine the transitions
 185 between states within the event. For example, in Figure 2A, the value of the Weather situation
 186 feature (sunny or rainy) determines which of the Mood states is visited (happy or angry). At each
 187 time point, the model observes the value of a randomly selected feature of the current situation
 188 and responds to a query about which state will be visited next. In the example shown in Figure
 189 2A, the agent first observes that Day = weekday, and then is asked to predict the upcoming Barista
 190 state (will the barista be Bob or Eve). Then it observes that Homework = paper and is asked to

191 predict the upcoming Mood state (will the barista be happy or angry). Finally, it observes that the
192 Weather = sunny and is asked to predict the upcoming Coffee state (will the drink be latte or mocha).
193 Both observations and queries are represented by one-hot vectors. In our simulations, the length
194 of the event graph is 16 and the number of states for each time point is 4. This means the number
195 of unique ways in which an event can unfold (depending on the features of the current situation)
196 is 4^{16} – far too many to memorize. As such, learning an effective representation of the event graph
197 (i.e., which states can occur at which time points, and how the state transitions depend on the
198 values of the situation features) is essential for predicting which state will come next. In our model,
199 this information is learned during the meta-training phase and stored in the cortical network's
200 weights (i.e., the model's semantic memory). As a terminological point, in this paper we use the
201 term *situation* to refer to the “ground truth” of the feature-value pairings for the current event, and
202 we use *situation model* to refer to the model's internal representation of the current situation in
203 working memory (i.e., in the LSTM cell state).

204 Figure 2B shows the way we simulated the three conditions from *Chen et al. (2016)*. In each
205 of the conditions, the agent processes three events. Importantly, for all of the conditions, we im-
206 posed (by hand) an encoding policy where the model stored an episodic memory (reflecting the
207 current contents of working memory – i.e., the LSTM cell state) on the final time point of each event.
208 This encoding policy was based on previous findings suggesting that episodic encoding takes place
209 selectively at the end of an event (*Ben-Yakov and Dudai, 2011; Ben-Yakov et al., 2013; Baldassano*
210 *et al., 2017; Ben-Yakov and Henson, 2018; Reagh et al., 2020*); we critically examine this assump-
211 tion in the *Benefits of selectively encoding at the end of an event* section below. In both the RM and
212 DM conditions, the agent first processes a distractor event (i.e., event a1), and then processes two
213 related events that are controlled by the same situation (i.e., event b1 and b2). These two related
214 events capture the two-part movie in the study by *Chen et al. (2016)*, in the sense that knowing
215 information from the first event (b1) will make the second event (b2) more predictable. Note that,
216 at the start of movie part 2 (b2), models in both the RM and DM conditions have access to a lure
217 episodic memory that was formed during the distractor event (a1), and also a target episodic mem-
218 ory that was formed during movie part 1 (b1). The main difference is that, in the DM condition, the
219 working memory state is flushed between part 1 and part 2 (by resetting the cell state of the LSTM),
220 whereas the flush does not occur in the RM condition; this flush in the DM condition is meant to
221 capture the effects of the one-day delay between parts one and two in the study by *Chen et al.*
222 *(2016)*. Finally, in the NM condition, the agent processes three events from three different situa-
223 tions. Therefore, during movie part 2, the agent has no information in working memory or episodic
224 memory pertaining to part 1. The model was trained (repeatedly) to predict upcoming states on all
225 three trial types before being tested on each of these trial types (see the *Model training and testing*
226 section in the *Methods*).

227 To summarize, the task environment used in our simulations captures how understanding of
228 naturalistic events and narratives depends on memory: It is necessary to remember observations
229 from the past (possibly from a large number of time points ago) in order to optimally predict the
230 future. For example, in the *Twilight Zone* episode used by *Chen et al. (2016)*, learning that the
231 servants are robots early in the episode helps the viewer predict how one character will react
232 when another character suggests killing all of the servants; similarly, in the model, learning that
233 the weather is sunny during event b1 will help the model predict that the barista will be happy
234 during event b2. The model is incentivized to routinely hold observations in working memory,
235 because information that is observed early in an event can sometimes be used to answer queries
236 that are posed later in that same event, or possibly across events (in the RM condition). This should
237 lead to a dynamic where the amount of information held in working memory builds within an
238 event (i.e., with each successive observation, the model builds a more “complete” representation in
239 working memory of the features of the current situation). Episodic memory is incentivized because
240 of the working memory “flush” in the DM condition between events b1 and b2 – information that is
241 relevant to b2 is observed during b1 but flushed from working memory, so the only way to benefit

242 from this information is to store it in episodic memory (at the end of b1) and then retrieve it from
243 episodic memory at the start of b2 (for additional discussion of how episodic memory can help to
244 bridge interruptions, see classic work by *Ericsson and Kintsch 1995*).

245 Figure 2C illustrates these points by showing the decoded contents of the model's working
246 memory for an example DM trial. To generate this figure, a linear classifier (logistic regression with
247 L2 penalty) was used to decode whether the correct (i.e., observed) value of each situation feature
248 was represented in the working memory state of the model (i.e., the LSTM cell state) at each time
249 point during the trial; see the *Decoding the working memory state* section in the *Methods* for more
250 details. We found that, once a feature was observed (indicated by a green box in the figure), this
251 feature typically was decodable until the end of the event, which confirms that observed features
252 tend to be actively maintained in the working memory state of the agent. The figure makes it clear
253 how, because of this tendency to maintain information over time, the model's representation of the
254 situation becomes more complete over time within part 1 of the event. The model then stores an
255 episodic memory snapshot on the final time point in part 1 (indicated by the blue arrow). Between
256 part 1 and part 2, the model's working memory state is flushed; then, early in part 2, the model
257 retrieves the stored episodic memory snapshot (indicated by the red arrow), which results in many
258 features of the situation becoming decodable before they are actually observed during part 2.

259 We acknowledge that our event-processing simulations incorporate several major simplifica-
260 tions. For example, we are modeling the first part of the movie as a single event when, in the *Chen*
261 *et al. (2016)* study, each half of the Twilight Zone episode clearly contains multiple events. We also
262 are assuming that the rate of key situation features being revealed is linear (one per time point)
263 and that feature values stay stable within events. Our goal here was to come up with the simplest
264 possible framework that allowed us to meaningfully engage with questions about encoding and
265 retrieval policies for episodic memory. In the *Discussion*, we talk about ways that the model could
266 be extended to more fully address the complexity of real-world events.

267 **The learned retrieval policy is sensitive to uncertainty**

268 Figure 3A shows the trained model's prediction performance during movie part 2, with the penalty
269 value for incorrect prediction set to 2. In the recent memory (RM) condition, prediction accuracy
270 is at ceiling starting from the beginning of part 2 – all situation feature values for the ongoing
271 situation were observed during the first part of the sequence, and the model is able to hold on
272 to these features in working memory. In the distant memory (DM) condition, prediction accuracy
273 starts out much lower, but after a few time points the accuracy is almost at ceiling. In the no
274 memory condition (NM), prediction accuracy increases linearly, reflecting the fact that the model
275 is able to observe more situation features as the event unfolds. The fact that prediction accuracy
276 is better in the DM condition than in the NM condition suggests that the model is using episodic
277 memory to support prediction in the DM condition.

278 We were particularly interested in whether the model's learned retrieval policy would be demand-
279 sensitive (i.e., would the model be more prone to retrieve from episodic memory if there were gaps
280 in its situation model, leading it to be uncertain about the upcoming state). To answer this ques-
281 tion, we visualized the activation levels of the target and lure memories during part 2, for each
282 of the three conditions (Figure 3B). Across the three conditions, we found much higher levels of
283 memory activation in the DM condition than the other two conditions. Importantly, the finding (in
284 the model) of greater memory activation in the DM condition than the RM condition qualitatively
285 captures the finding from *Chen et al. (2016)* that the putative fMRI signature of episodic retrieval
286 (hippocampal-cortical coupling) was stronger in the DM condition than the RM condition. Note that,
287 in our simulation, the set of available episodic memories in the RM and the DM condition is the
288 same. The main difference is that, in the RM condition, the network has a fully-specified situation
289 model actively maintained in its working memory (the recurrent activity of the LSTM) during part
290 2, which is sufficient for the network to predict the upcoming state. In contrast, at the beginning
291 of the DM condition, the network's ongoing situation model is empty – the values for all features

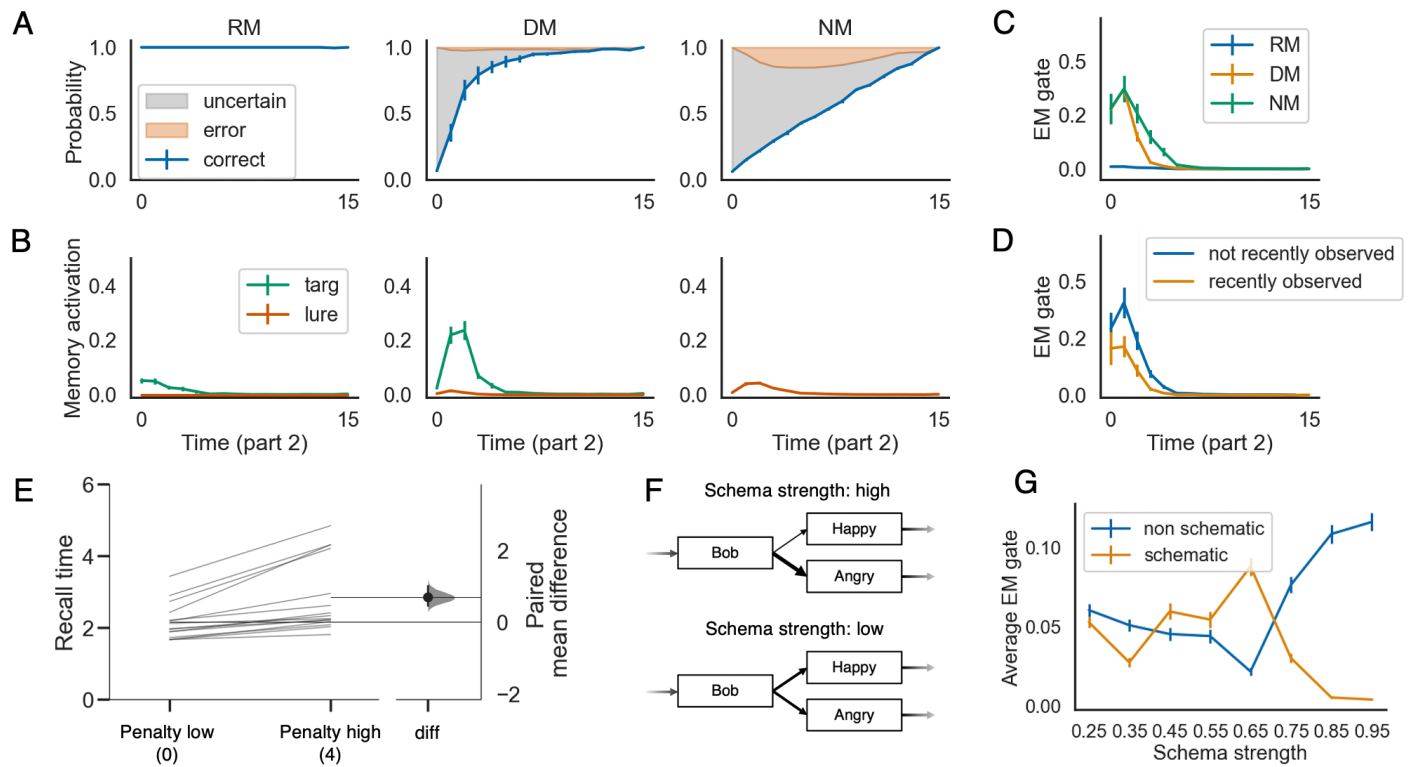


Figure 3. The learned episodic retrieval policy is selective. Panels A, B, and C show the model's behavioral performance, memory activation, and episodic memory gate (EM gate) value during part 2, across the recent memory (RM), distant memory (DM), and no memory (NM) conditions, when the penalty for incorrect prediction is set to 2 at test. These results show that recall is much stronger in the DM condition (where episodic retrieval is needed to fill in gaps and resolve uncertainty) compared to the RM condition. D) shows that, in the DM condition, the EM gate value is lower if the model has recently (i.e., in the current event) observed the feature that controls the upcoming state transition. E) shows how the average recall time is delayed when the penalty for making incorrect predictions is higher. F) illustrates the definition of the schema strength for a given time point. G) shows how the average EM gate value changes as a function of schema strength (penalty level = 2). The errorbars indicate 1SE across 15 models.

292 are unknown. Overall, this result suggests that the learned retrieval policy is demand-sensitive (for
293 simulations of other, related findings from this study, see *Appendix 6*).

294 To gain further insight into the model's retrieval policy, we examined the EM gate values in
295 the three conditions (Figure 3C). We found that the model sets the EM gate to a higher (more
296 "open") value in the DM and NM conditions (where there are gaps in the model's understanding of
297 the ongoing situation, causing it to be uncertain about what was coming next), and it suppresses
298 episodic retrieval in the RM condition (where there are no gaps). Likewise, within the DM condi-
299 tion, the model sets the EM gate to a higher value when the feature controlling the next transition
300 has not been recently observed (i.e., the feature is not in working memory, causing the model to
301 be uncertain about what was coming next) vs. if the relevant feature has been recently observed
302 and is therefore active in working memory (Figure 3D). The same principle also explains why, for
303 later time points in part 2, the EM gate is set to a lower value in the DM condition than the NM
304 condition (Figure 3C) – in the DM condition, episodic retrieval that occurs on earlier time points
305 makes the model more certain on later time points, reducing the demand for episodic retrieval
306 and (consequently) leading to lower EM gate values.

307 The fact that the model learned a demand-sensitive retrieval policy can be explained in terms
308 of a simple cost-benefit analysis: When the model is unsure about what will happen next, the
309 potential benefits of episodic retrieval are high. In the absence of episodic retrieval, the model
310 will have to guess or say "don't know", but if it consults episodic memory, the model could end up
311 recalling the feature of the situation that controls the upcoming state transition, allowing it to make
312 a correct prediction. By contrast, when the feature of the situation that controls the transition is
313 already in working memory (and consequently the model is able to make a specific prediction about
314 what will happen next), there is less of a benefit associated with episodic retrieval – the only way
315 that episodic retrieval will help is if the model is holding the wrong feature in working memory and
316 the episodic memory overwrites it. Furthermore, in this scenario, there is also a potential cost to
317 retrieving from episodic memory: Lures are always present, and if the model recalls a lure this can
318 overwrite the correct information in working memory. Since the potential costs of episodic retrieval
319 outweigh the benefits of episodic retrieval in the "high certainty" scenario, the model learns a policy
320 of waiting to retrieve until it is uncertain about what will happen next.

321 Importantly, the model's ability to *adjust its policy* when it is uncertain is predicated on there
322 being a reliable "neural correlate of certainty" in the model, which can be used as the basis for this
323 differential responding; we investigated this and found that the norm of activity in the decision
324 layer is lower when the model is uncertain vs. certain (for more details, see *Appendix 1*). This (im-
325 plicit) neural correlate of certainty exists regardless of whether the model is trained to explicitly
326 signal uncertainty via the "don't know" response. In other simulations (reported in *Appendix 5*), we
327 found that a version of the model without the "don't know" option can still leverage this implicit
328 neural correlate of certainty to show demand-sensitive retrieval (i.e., more episodic retrieval in
329 the DM condition than the RM condition); the main effect of including the "don't know" option is
330 to make the model more patient overall, by reducing the cost associated with waiting to retrieve
331 from episodic memory.

332 **The effect of penalty on retrieval policy**

333 A key question is how the model's policy for prediction and episodic retrieval adapts to different
334 environmental regimes. Toward this end, we explored what happens when we vary the penalty on
335 false recall from 0 to 4 during model meta-testing – that is, can the model flexibly adjust its policy
336 based on the current penalty? (note that the penalty was uniformly sampled from the 0-4 range
337 during meta-training). If learning a selective retrieval policy is driven by the need to manage the
338 costs of false recall, then it stands to reason that varying these costs should affect the model's policy.
339 Our first finding is that adjusting the penalty at test affects the model's tendency to give "don't
340 know" responses: When the penalty is zero, the model makes specific next-state predictions (i.e., it
341 refrains from using the "don't know" response) even when it can not reliably predict the next state,

342 leading to many errors. In contrast, when the penalty is high, the model makes more “don’t know”
343 responses (in the DM condition, the model responds “don’t know” 15.8% of the time when penalty
344 is set to 4, vs. 0.3% of the time when penalty is set to 0). This strategy is rational – when the penalty
345 is zero, the expected reward is larger for randomly guessing an answer than for saying “don’t know”,
346 but when the penalty is set to four, the expected reward is larger for saying “don’t know” than for
347 random guessing. We also found that, when the model is tested in an environment where the
348 penalty is high, it waits longer to retrieve from episodic memory, relative to when the penalty at
349 training is lower (Figure 3E). This delay in recall can be explained in terms of a speed-accuracy
350 trade-off. Waiting longer to retrieve from episodic memory allows the model to observe more
351 features, which helps to disambiguate the present situation from other, related situations and
352 thereby reduces false recall. However, waiting longer also carries an opportunity cost – the model
353 has to forego all of the rewards it would have received (from correct prediction) if it had recalled
354 the correct memory earlier. When the penalty is low, the benefits of retrieving early (in terms of
355 increased correct prediction) outweigh the costs (in terms of increased incorrect prediction due to
356 false recall), but when the penalty is high, the costs outweigh the benefits, so the model is more
357 cautious and it waits to observe more features to be sure that the memory it (eventually) recalls is
358 the right one.

359 **The effect of schema regularity on the learned policy**

360 Next, we examined the effect of schema regularity on the agent’s retrieval policy. In the simulations
361 preceding this one, we imposed a form of schematic structure by teaching the model about which
362 states could be visited at which time points (i.e., the “columns” of Figure 2A). However, *within* a
363 particular time point, the marginal probabilities of the states that were “allowed” at that time point
364 were equated – put another way, none of the states were more prototypical than any of the other
365 states. In this simulation, we also allowed for some states to be more prototypical (i.e., occur more
366 often) than other states that could occur at that time point. We say that a time point is *schematic*
367 if there is one state that happens with higher probability, compared to other states. Consider the
368 example illustrated in Figure 3F: If the probability of Bob being angry is much greater than the
369 probability of him being happy, then we say that this is a highly schematic time point. In contrast,
370 if Bob is equally likely to be happy or angry, then the schema strength is low. Intuitively, when
371 there is a strong schema, there is less of a need to rely on episodic memory – in the limiting case, if
372 the schematic state occurs in every sequence, the model will learn to predict this state every time
373 and there is no need to consult episodic memory.

374 To explore the effects of schema strength, we ran simulations where half of the time points
375 were schematic. For the other half of the time points (*non-schematic* time points), all of the states
376 associated with that time point were equally probable (given that there were four possible states at
377 each time point, the probability of each state was .25). Schematic and non-schematic time points
378 were arranged in an alternating fashion (for half of the models, even time points were schematic
379 and odd time points were non-schematic, and the opposite was true for the other half of the mod-
380 els). For schematic time points, we manipulated the strength of schematic regularity in the envi-
381 ronment by manipulating the probability of the “prototypical” state. We varied schema strength
382 values from 0.25 (baseline) to 0.95 in steps of 0.10.

383 The results of this analysis when penalty was set to 2 at test are shown in Figure 3G, which
384 plots the EM gate value during part 2 as a function of schema strength. The first thing to note
385 about these results is that, for high levels of schema strength, episodic retrieval is suppressed for
386 schematic time points (i.e., time points with a prototypical state) and elevated for non-schematic
387 time points (i.e., time points where there was not a prototypical state). The former finding (sup-
388 pression of retrieval at time points where there is a strong prototype) fits with the intuition, noted
389 above, that high-schema-strength states are almost fully predictable without episodic memory,
390 and thus there is no need to retrieve from episodic memory. The latter finding (enhanced retrieval
391 at non-schematic time points, when schema strength is high overall) can be explained in terms

392 of the idea that schema-congruent features tend to be shared by both target and lure memories
393 and thus are not diagnostic of which memory is the target; in this situation, the only way to distin-
394 guish between targets and lures is to recall non-schematic features, which is why the model tries
395 extra-hard to retrieve them from episodic memory.

396 Interestingly, the model shows the opposite pattern of effects when schema strength = .55 or
397 .65: Episodic retrieval is enhanced for schematic time points and suppressed for non-schematic
398 time points. This reversal can be explained as follows: When schema strength = .55 or .65, the
399 model has started to build up a tendency to guess the schema-congruent (prototypical) state, but
400 it is also going to be wrong about 1/3 of the time when it guesses the schema-congruent state,
401 incurring a substantial penalty. To counteract this tendency to make wrong guesses, the model
402 needs to try extra-hard to retrieve the actual feature value for schematic time points (which is why
403 the EM gate value increases for these time points) – and if the model is doing more retrieval in
404 response to schematic states, it needs to do somewhat less retrieval in response to non-schematic
405 states (which is why the EM gate value goes down for these features). As schema strength increases
406 beyond .65, the model will be wrong less often when it guesses the schema-congruent state, so
407 there is less of a need to counteract wrong guesses with episodic retrieval – this makes it safe
408 for the model to reduce the EM gate value for schematic time points at higher levels of schema
409 strength (as described above).

410 **Other factors that affect the learned retrieval policy**

411 In addition to the simulations described above, we also ran simulations exploring the effects of
412 *between-event similarity* and *familiarity* on the learned retrieval policy. With regard to similarity: We
413 found that the model is more cautious about retrieving from episodic memory if trained in environ-
414 ments where memories are highly similar (because the risk of false recall is higher) – see *Appendix*
415 *2* for details. With regard to familiarity: When we provided the model with a familiarity signal that
416 is informative about whether a situation was previously encountered, we found that the model
417 learns to exploit this information by retrieving more from episodic memory when the familiarity
418 signal is high and retrieving less from episodic memory when the familiarity signal is low. This
419 result provides a resource-rational account of experimental findings showing that familiar stimuli
420 shift the hippocampus into a “retrieval mode” where it is more likely to (subsequently) retrieve
421 episodic memories (*Duncan et al., 2012; Duncan and Shohamy, 2016; Duncan et al., 2019; Patil*
422 *and Duncan, 2018; Hasselmo and Wyble, 1997*) – see *Appendix 3* for details.

423 **Benefits of selective encoding**

424 Above, we showed that the model learned selective retrieval policies (e.g., avoiding retrieval from
425 episodic memory early on during part 2, or when certain about upcoming states) in order to reduce
426 the risk of recalling irrelevant memories. Here, we shift our focus to the complementary question
427 of *encoding policy*: When is the best time to store episodic memories? In the simulations reported
428 below, we show that a selective encoding policy can benefit performance, by reducing interference
429 at retrieval later on. Note that our model is presently not capable of learning an encoding policy on
430 its own (see *Discussion*), but we can explore the benefits of selective encoding by imposing different
431 encoding policies by hand and seeing how they affect performance.

432 **Benefits of selectively encoding at the end of an event**

433 The simulations presented thus far assumed that episodic memories are selectively encoded at
434 the ends of events. This assumption was based on findings from several recent fMRI studies that
435 measured hippocampal activity during perception of events and related this to later memory for
436 the events. These studies found that the hippocampal response tends to peak at event bound-
437 aries (*Ben-Yakov and Dudai, 2011; Ben-Yakov et al., 2013; Baldassano et al., 2017; Ben-Yakov and*
438 *Henson, 2018; Reagh et al., 2020*); this boundary-locked response predicts subsequent memory
439 performance for the just-completed event (*Ben-Yakov and Dudai, 2011; Baldassano et al., 2017;*

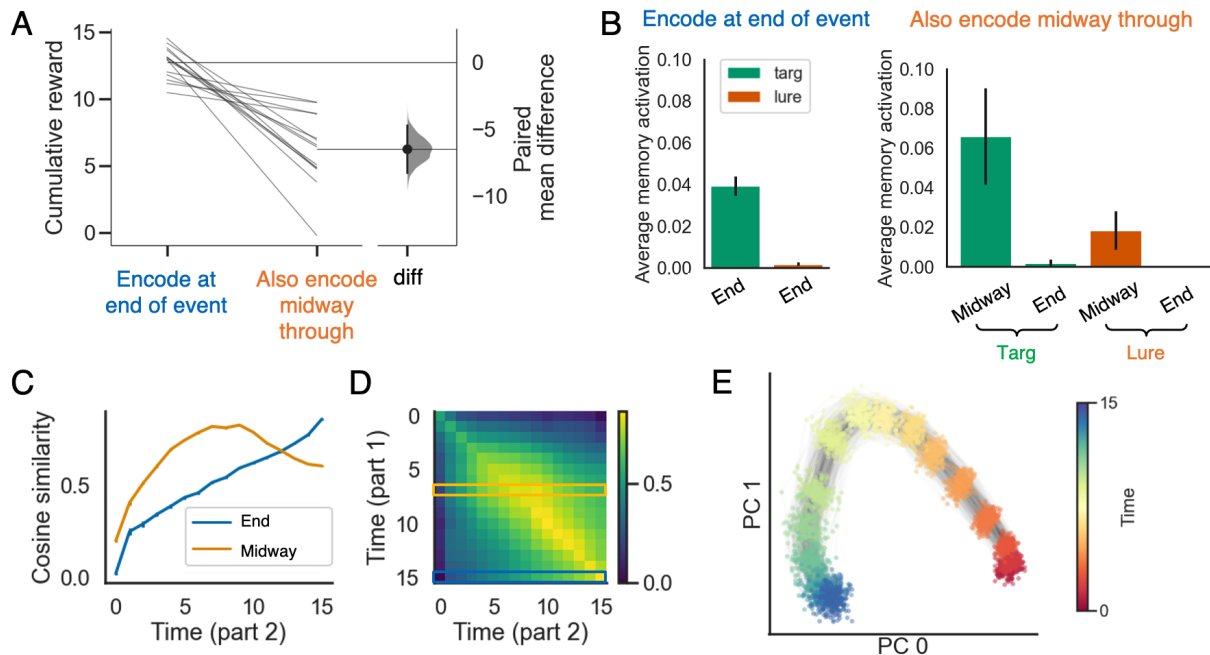


Figure 4. The advantage of selectively encoding episodic memories at the end of an event. A) Prediction performance is better for models that selectively encode at the end of each event, compared to models that encode at the end of each event and also midway through each event. B) The model performs worse with midway-encoded memories because midway-encoded target memories are activated more strongly than end-encoded target memories, thereby blocking recall of the (more informative) end-encoded target memories, and also because midway-encoded lure memories are more strongly activated than end-encoded lure memories (see text for additional discussion). C) The cosine similarity between working memory states during part 2 and memories formed midway through part 1 (in orange) or at the end of part 1 (in blue). The result indicates that the midway-encoded memory will dominate the end-encoded memory for most time points. D) The time-point-to-time-point cosine similarity matrix between working memory states from part 1 versus part 2 in the no memory (NM) condition (part C depicts the orange and blue rows from this matrix). E) PCA plot of working memory states as a function of time, for a large number of events. The plot shows that differences in time within events are represented much more strongly than differences across events. The errorbars indicate 1SE across 15 models.

440 *Reagh et al., 2020*), leading researchers to conclude that it is a neural signature of episodic encod-
441 ing of the just-completed event.

442 While these results suggest that the end of an event may be a particularly important time for
443 episodic encoding, existing studies do not provide a computational account of *why* this should
444 be the case. This “why” question can be broken into two parts: First, why might it be beneficial
445 to encode at the end of an event, and second, why might it be *harmful* to encode at other times
446 within the event? Answering the first question (regarding benefits of encoding at the end of an
447 event) is relatively straightforward. Several researchers have argued that information about the
448 current situation builds up in working memory within an event, and then is “flushed out” at event
449 boundaries (*Radvansky et al. 2011; Richmond and Zacks 2017*; for neural evidence in support of
450 this dynamic, see *Ezzyat and Davachi 2011; Chien and Honey 2020; Ezzyat and Davachi 2021*). This
451 dynamic (which is illustrated in the model in Figure 2C) means that the model’s representation of
452 the features of an event will be most complete right before the end of the event, making this a
453 particularly advantageous time to take an episodic memory snapshot of the situation model.

454 While it is clear why encoding at the end of an event is useful, it is less clear why encoding
455 at other times might be harmful; naively, one might think that storing more episodic snapshots
456 during an event would lead to *better* memory for the event. To answer this question, we compared
457 models that selectively encode episodic memories at the end of each event to models that encode
458 episodic memories both at the end of each event and also midway through each event. If selectively
459 encoding at the end of an event yields better performance, this would provide a resource-rational
460 justification for the empirical findings reviewed above.

461 Our simulation shows that in the DM condition, during part 2, models that encode an addi-
462 tional episodic memory midway through each event performed worse (Figure 4A). This decrease
463 in performance can be explained in terms of several related factors. First, as shown in Figure 4B,
464 when midway memories are also stored, midway memories of the target event are recalled more
465 strongly than memories formed at the end of the target event.

466 This advantage occurs because the model’s hidden state strongly encodes temporal context:
467 WM states stored at similar times within an event tend to be more similar than WM states stored
468 at different times (this illustrated by Figure 4E, which shows that time information within an event is
469 more strongly represented than differences across events). This strong temporal encoding makes
470 sense, given that the model needs to know where it is in the sequence in order to predict which
471 observations will come next (for a review of evidence for this kind of temporal coding in the brain,
472 see *Eichenbaum 2014*). One consequence of this time coding is that – early on in part 2 of the
473 event (when the benefits of episodic retrieval are the largest) – the temporal context represented
474 in working memory will be a better match for memories encoded midway through the event than
475 memories encoded at the end of the event (Figure 4C and D). This temporal context match pro-
476 vides a competitive advantage for the midway memory over the endpoint memory, resulting in
477 the midway memory blocking the endpoint memory from coming strongly to mind. The second
478 key point is that the midway memory is less informative (i.e., it contains fewer features of the sit-
479 uation, because it was stored before the full set of features was observed). As such, recalling the
480 midway target memory confers less of a benefit on future prediction than recalling the endpoint
481 memory would have provided – this is the main reason why prediction is worse in the midway con-
482 dition. The third key point is that, because midway memories contain less information, they are
483 more confusable across events (i.e., it is harder to determine which event the memory pertains to).
484 As a result, midway lures tend to become more active at retrieval than endpoint lures (Figure 4B) –
485 this lure retrieval acts to further reduce prediction accuracy.

486 A possible alternative explanation of the negative effects of midway encoding is that midway
487 encoding was introduced when we tested the model’s performance but was not present during
488 meta-training (i.e., when the model acquired its retrieval policy); as such, midway encoding can be
489 viewed as “out of distribution” and may be harmful for this reason. To address this concern, we also
490 ran a version of the model where memories were stored both midway and at the end of an event

491 during meta-training, and it was still true that endpoint-only encoding led to better performance
492 than midway-plus-endpoint encoding; this result shows that midway encoding is intrinsically harm-
493 ful, and it is not just a matter of it being out-of-distribution.

494 To summarize the results from this simulation, the model does better when we force it to wait
495 until the end of an event to take a snapshot; this occurs because midway target memories block
496 recall of more informative endpoint target memories, and also because there is more false recall
497 of midway lures than endpoint lures. This model result provides a resource-rational justification
498 for the results cited above showing preferential encoding-related hippocampal activity at the end
499 of events (*Ben-Yakov and Dudai, 2011; Ben-Yakov et al., 2013; Baldassano et al., 2017; Ben-Yakov*
500 *and Henson, 2018; Reagh et al., 2020*).

501 Discussion

502 Most of what we know about episodic memory has, by design, come from experiments where
503 performance depends primarily on episodic memory (as opposed to other memory systems), and
504 participants are given clear instructions about when episodic memories should be stored and re-
505 trieved (e.g., learning and recalling lists of random word pairs); likewise, most computational mod-
506 els of human memory have focused on explaining findings from these kinds of experiments (for a
507 review, see *Norman et al. 2008*). However, as noted in the *Introduction*, real-world memory does
508 not adhere to these constraints: In naturalistic learning situations, participants are typically not
509 given any instructions about how episodic memory should be used to support performance, and
510 – even when participants are given instructions about what to remember – performance usually
511 depends on a complex mix of memory systems, with contributions from both working memory
512 and semantic memory in addition to episodic memory.

513 The goal of the present work was to gain some theoretical traction on when episodic memo-
514 ries should be stored and retrieved to optimize performance in these more complex situations.
515 Towards this end, we optimized a neural network model that *learned its own policy* for when to con-
516 sult episodic memory (via an adjustable gate) in order to maximize reward, and we also (by hand)
517 explored the effects of different episodic memory encoding policies on network performance. Our
518 approach is built on the principle of resource rationality, whereby human cognition is viewed as
519 an approximately-optimal solution to the learning challenges posed by the environment, subject
520 to constraints imposed by our cognitive architecture (*Griffiths et al., 2015; Lieder and Griffiths,*
521 *2019*); according to this principle, the approximately-optimal solutions obtained by our model can
522 be viewed as hypotheses about (and explanations of) how humans use episodic memory in com-
523 plex, real-world tasks.

524 In the simulations presented here, we identified several ways in which selective policies for
525 episodic memory retrieval and encoding can benefit performance. With regard to retrieval, we
526 showed that the model learns to avoid episodic retrieval in situations where the risks of retrieval
527 (i.e., retrieving the wrong memory, leading to incorrect predictions) outweigh the benefits (i.e., re-
528 trieving the correct memory, leading to increased correct predictions). For example, when there
529 is high certainty about what will be observed next (due to the relevant information being main-
530 tained in working memory or semantic memory), the marginal benefits of retrieving from episodic
531 memory are too small to outweigh the risks of retrieving the wrong memory. Another example
532 is when too little information has been observed to pinpoint the relevant memory – in this case,
533 the potential benefits of retrieving are high, but the risks of retrieving the wrong memory are also
534 high, leading the model to defer retrieving until more information has been observed. With re-
535 gard to encoding, we showed that waiting until the end of an event to encode a memory for that
536 event boosts subsequent prediction performance – this performance boost comes from reducing
537 “clutter” (interference) from other memories, thereby making it easier to retrieve the sought-after
538 memory. These modeling results explain a wide range of existing behavioral and neuroimaging re-
539 sults, and also lead to new, testable predictions. With regard to existing results: The model provides
540 a resource-rational account of findings from *Chen et al. (2016)* showing the demand-sensitivity of

541 episodic retrieval, as well as results showing that episodic encoding is modulated by event bound-
542 aries (*Ben-Yakov and Dudai, 2011; Ben-Yakov et al., 2013; Baldassano et al., 2017; Ben-Yakov and*
543 *Henson, 2018; Reagh et al., 2020*). Appendix 3 also shows how the model explains effects of famil-
544 iarity on retrieval policy (*Duncan et al., 2012; Duncan and Shohamy, 2016; Duncan et al., 2019; Patil*
545 *and Duncan, 2018; Hasselmo and Wyble, 1997*). With regard to novel predictions: Our model makes
546 predictions about how episodic retrieval will be modulated by certainty (Figure 3B, C, D), penalty
547 (Figure 3E), schema strength (Figure 3G), and similarity (Figure 1) – all of these predicted relation-
548 ships could be tested in experiments that measure hippocampal-cortical information transfer, ei-
549 ther using measures like hippocampal-cortical inter-subject functional connectivity in fMRI (e.g.,
550 *Chen et al. 2016; Chang et al. 2021*) or time-lagged mutual information in ECoG (e.g., *Michelmann*
551 *et al. 2021*).

552 More broadly, the simulations presented here show how the model can be used to explore in-
553 teractions between three distinct memory systems: semantic memory (instantiated in the weights
554 in cortex), working memory (instantiated in the gating policy learned by the cortical LSTM module,
555 allowing for activation at one time point in cortex to influence activation at subsequent time points),
556 and episodic memory. In the past, modelers have focused on these memory systems in isolation
557 (see, e.g., *Norman et al. 2008*), in part because of a desire to understand the detailed workings
558 of the systems, but also because of technical limitations: Until very recently, the technology did
559 not exist to automatically optimize the performance of networks containing episodic memory, so
560 researchers interested in simulating interactions between episodic memory and these other sys-
561 tems were put in the position of having to do time-consuming (and frustrating) hand-optimization
562 of the models. Here, we leverage recent progress in the artificial intelligence literature on memory-
563 augmented neural networks (*Graves et al., 2016; Pritzel et al., 2017; Ritter et al., 2018; Wayne*
564 *et al., 2018*) that makes it possible to automatically optimize the use of episodic memory and its
565 interactions with other memory systems. This technical advance has opened up a new frontier in
566 the cognitive modeling of memory (*Collins, 2019*), making it possible to address both “naturalistic
567 memory” scenarios and controlled experiments that involve interactions between prior knowledge
568 (semantic memory), active maintenance (working memory), and episodic memory.

569 **Relation to other models**

570 Memory-augmented neural networks with a differentiable neural dictionary
571 Conceptually, the episodic memory system used in our model is similar to recently-described
572 memory-augmented neural networks with a differentiable neural dictionary (DND) (*Pritzel et al.,*
573 *2017; Ritter et al., 2018; Ritter, 2019*). In these models, the data structure of the episodic memory
574 system is dictionary-like: Each memory is a key-value pair. The keys define the similarity metric
575 across all memories, and the values represent the content of these memories. For example, one
576 can use the LSTM cell state patterns as the keys and use the final output of the network as the
577 values (*Pritzel et al., 2017*); note that, in our model, the cell state of the cortical network serves as
578 both the key and the value. The work by *Ritter et al. (2018)* is particularly relevant as it was the first
579 paper (to our knowledge) to use the DND for cognitive modeling and – as such – served as a major
580 inspiration for the work presented here (see also *Botvinick et al. 2019*). The way that our model
581 uses the DND mechanism is quite similar to how it was used in *Ritter et al. (2018)*; in particular,
582 we took from the *Ritter et al. (2018)* paper the idea that the cortical network learns to control a
583 “gate” on episodic retrieval via reinforcement learning. However, there are also some meaningful
584 differences between our model and the model used by *Ritter et al. (2018)*.

585 The most salient difference regards the placement of the EM gate: In our model, the gate con-
586 trols the flow of information into the episodic memory module (*pre-gating*), but in the Ritter model
587 the gate controls the flow of information *out* of the episodic memory module (*post-gating*). Practi-
588 cally speaking, the main consequence of having the gate on the output side is that the gate can be
589 controlled based on information coming out of the hippocampus, in addition to all of the cortical
590 regions that are used to control the gate in our pre-gating model. While this is a major difference,

591 we found that our key simulation results qualitatively replicate in a version of the model that uses
592 post-gating, indicating that the selective encoding and retrieval principles discussed here do not
593 depend on the exact placement of the gate (see *Appendix 5* for simulation results and more discus-
594 sion of these points).

595 Another difference is that our model's computation of which memories are retrieved (given a
596 particular retrieval cue, assuming that the "gate" on retrieval is open) is more complex. *Ritter et al.*
597 (2018) used a one-nearest-neighbor matching algorithm during recall, whereby the stored memory
598 with the highest match to the cue is selected for retrieval (assuming that the gate is open). By
599 contrast, memory activation in our model is computed using a competitive evidence accumulation
600 process, in line with prior cognitive models of retrieval (e.g., *Polyn et al. 2009; Sederberg et al.*
601 *2008*). While we did not explore the effects of varying the level of competition in our simulations,
602 having this as an adjustable parameter opens the door to future work where the model learns a
603 policy for setting competition in order to optimize performance (just as it presently learns a policy
604 for setting the EM gate).

605 A third structural difference between our model and the *Ritter et al. (2018)* model is our addi-
606 tion of the "don't know" output unit, which (when selected) allows the model to avoid both reward
607 and punishment. As discussed above, the primary effect of incorporating this "don't know" action
608 is to make the model more patient (i.e., more likely to wait to retrieve from episodic memory), by
609 giving it a way to avoid incurring penalties if it decides to wait to retrieve (for more details, see
610 *Appendix 5*).

611 Apart from the structural differences noted above, the main difference between our model-
612 ing work and the work done by *Ritter et al. (2018)* relates to the application domain (i.e., which
613 cognitive phenomena were simulated). Our modeling work in this paper focused on how episodic
614 memory can support incidental prediction of upcoming states, when there is no explicit demand
615 for a decision. By contrast, *Ritter et al. (2018)* focused on how episodic memory can be used to
616 support performance in classic decision-making tasks, such as bandit tasks and maze learning, that
617 have been extensively explored in the reinforcement learning literature.

618 The structured event memory (SEM) model

619 Another highly relevant model is the structured event memory (SEM) model developed by *Franklin*
620 *et al. (2020)*. Like our model, SEM uses RNNs to represent its knowledge of schemas (i.e., how
621 events typically unfold). Also, like our model, SEM records episodic memory traces as it processes
622 events. However, there are several key differences between our model and SEM. First, whereas
623 our model uses a single RNN to represent a single (contextually parameterized) schema, SEM uses
624 multiple RNNs that each represent a distinct schema for how events can unfold. Building on prior
625 work on nonparametric Bayesian inference (*Anderson, 1991; Aldous, 1985; Pitman, 2006*) and latent
626 cause modeling (*Gershman et al., 2010, 2015*), SEM contains specialized computational machinery
627 that allows it to determine which of its stored schemas (each with its own RNN) is relevant at a
628 particular moment, and also when it is appropriate to instantiate a new schema (with its own, new
629 RNN) to learn about ongoing events. This inference machinery allows SEM to infer when event
630 boundaries (i.e., switches in the relevant schema) have occurred; the *Franklin et al. (2020)* paper
631 leverages this to account for data on how people segment events. Our model lacks this inference
632 machinery, so we need to impose event boundaries by fiat, as opposed to having the model identify
633 them on their own.

634 Another major difference between the models relates to how episodic memory is used. A key
635 focus of our modeling work in this paper is on how episodic memory can support online prediction.
636 By contrast, in SEM, episodic memory is not used at all for online prediction – online prediction is
637 based purely on the weights of the RNNs (i.e., semantic memory) and the activation patterns in the
638 RNNs (i.e., working memory). The sole use of episodic memory in the *Franklin et al. (2020)* paper
639 is to support reconstruction of previously-experienced events. Specifically, in SEM, each time point
640 leaves behind a noisy episodic trace; the *Franklin et al. (2020)* paper shows how Bayesian inference

641 can combine these noisy stored episodic memory traces with stored knowledge about how events
642 typically unfold (in the RNNs) to reconstruct an event. Effectively, SEM uses knowledge in the RNNs
643 to “de-noise” and fill in gaps in the stored episodic traces. The *Franklin et al. (2020)* paper uses this
644 process to account for several findings relating to human reconstructive memory.

645 **Future directions and limitations**

646 On the modeling side, our work can be extended in several different ways. As noted above, our
647 model and SEM have complementary strengths: SEM is capable of storing multiple schemas and
648 doing event segmentation, whereas our model only stores a single schema and we impose event
649 boundaries by hand; our model is capable of using episodic memory to support online prediction,
650 whereas SEM is not. It is easy to see how these complementary strengths could be combined into
651 a single model: By adding SEM’s ability to do multi-schema inference to our model, we would be
652 able to simulate both event segmentation and the role of episodic memory in predicting upcoming
653 states, and we would also be able to explore *interactions* between these processes (e.g., using
654 episodic memory to predict could affect when prediction errors occur, which – in turn – could
655 affect how events are segmented; *Zacks et al. 2007, 2011*).

656 Another limitation of the current model is that the encoding policy is not learned. In our sim-
657 ulations, we trained models with different (pre-specified) encoding policies and compared their
658 performance. Going forward, we would like to develop models that learn when to encode through
659 experience, instead of imposing encoding policies by hand. Our results show that selective encod-
660 ing can yield better performance than encoding everything, so – in principle – selective encoding
661 policies should be learnable with RL. The main challenge in learning encoding policies is the long
662 temporal gap between the decision to encode (or not) and learning the consequences of that choice
663 for retrieval. Moreover, a high-quality encoding policy, taken on its own, generally does not lead to
664 high reward when the retrieval policy is bad; that is, encoding policy and retrieval policy have to be
665 learned in a highly coordinated fashion. Recent technical advances in RL (e.g., algorithms that do
666 credit assignment across long temporal gaps; *Raposo et al. 2021*) may make it easier to address
667 these challenges going forward.

668 A benefit of being able to learn encoding policies in response to different task demands is that
669 the model could discover other factors that it could use to modulate encoding – for example, sur-
670 prise. Numerous studies have found improved memory for surprising events (e.g., *Greve et al.*
671 *2017, 2019; Quent et al. 2021a; Kafkas and Montaldi 2018; Frank et al. 2020; Rouhani et al. 2018,*
672 *2020; Chen et al. 2015a; Pine et al. 2018; Antony et al. 2021; for reviews, see Frank and Kafkas*
673 *2021; Quent et al. 2021b*) – these behavioral results converge with a large body of literature show-
674 ing increased hippocampal engagement in response to prediction error (e.g., *Axmacher et al. 2010;*
675 *Chen et al. 2015a; Long et al. 2016; Kumaran and Maguire 2007, 2006, 2007; Duncan et al. 2012;*
676 *Davidow et al. 2016; Kafkas and Montaldi 2015; Frank et al. 2021; for reviews, see Frank and*
677 *Kafkas 2021; Quent et al. 2021b*), and also with a recent fMRI study showing that prediction error
678 biases hippocampal dynamics towards encoding (*Bein et al., 2020*). Given that studies have found
679 a strong relationship between surprise and event segmentation (e.g., *Zacks et al. 2007, 2011; for a*
680 *recent example see Antony et al. 2021*), it seems possible that increased episodic encoding at the
681 ends of events could be driven by peaks in surprise that occur at event boundaries. However, there
682 are complications to this view; in particular, some recent work has argued that not all event bound-
683 aries are surprising (*Schapiro et al., 2013*) – in light of this, more research is needed to explore the
684 relationship between these effects.

685 In addition to surprise, recent work by *Sherman and Turk-Browne (2020)* suggests that *predictive*
686 *certainty* may play a role in shaping encoding policy: They found that stimuli that trigger strong
687 predictions (i.e., high certainty about upcoming events) are encoded less well. In keeping with this
688 point, *Bonasia et al. (2018)* found that, during episodic encoding, events that were more typical
689 (and thus were associated with more predictive certainty, and less surprise) were associated with
690 lower levels of medial temporal lobe (MTL) activation. Intuitively, it makes sense to focus episodic

691 encoding on time periods where there is high surprise and low predictive certainty – if events in a
692 sequence are unsurprising and associated with high predictive certainty, this means that existing
693 (cortical) schemas are sufficient to reconstruct that event, and no new learning is necessary (or, if
694 learning is required, it is possible that cortex could handle this “schema-consistent” learning on its
695 own; *McClelland 2013; McClelland et al. 2020*). Conversely, if events in a sequence do not follow a
696 schema (leading to uncertainty) or violate that schema (leading to surprise), the only way to predict
697 those events later will be to store them in episodic memory. Future work can explore whether a
698 model that represents surprise and certainty (either implicitly or explicitly) can learn to leverage
699 one or both of these factors when deciding when to encode; our present model is a good place to
700 start in this regard, as we have already demonstrated the model’s ability to factor certainty into its
701 retrieval policy.

702 Another major simplification in the model’s encoding policy is that it stores each episodic mem-
703 ory as a distinct entity (see Figure 1B). Old memories are never overwritten or updated. However,
704 a growing literature on memory reconsolidation suggests that memory reminders can result in
705 participants accessing an existing memory and then updating that memory, rather than forming
706 a new memory outright (*Dudai and Eisenberg, 2004; Dudai, 2009; Hardt et al., 2010; Wang and*
707 *Morris, 2010*). In the future, we would like to develop models that decide whether to encode a
708 new episodic memory (pattern separate) or update an old memory (pattern complete). We could
709 implement this by having the model try to retrieve before it encodes a new memory; if it succeeds
710 in retrieving a stored memory above a certain threshold level of activation, the model could up-
711 date that memory rather than creating a new memory. In future work, we plan to implement this
712 mechanism and use it to simulate memory reconsolidation data.

713 Going forward, we also hope to explore more biologically-realistic episodic memory models
714 (e.g., *Schapiro et al. 2017; Norman and O’Reilly 2003; Ketz et al. 2013*). Using a more biologically-
715 realistic hippocampus could affect the model’s predictions (e.g., if memory traces were allowed to
716 interfere with each other during storage – currently they only interfere at retrieval) and it would
717 also improve our ability to connect the model to neural data on hippocampal codes and how they
718 change with learning (e.g., *Duncan and Schlichting 2018; Brunec et al. 2020; Ritvo et al. 2019; Fav-*
719 *ila et al. 2016; Chanales et al. 2017; Schlichting et al. 2015; Whittington et al. 2020; Stachenfeld*
720 *et al. 2017; Hulbert and Norman 2015; Kim et al. 2017; Schapiro et al. 2016, 2012*). Similarly, using
721 a more biologically-detailed cortical model (separated into distinct cortical sub-regions) could help
722 us to connect to data on how different cortical regions interact with hippocampus during event pro-
723 cessing (e.g., *Ranganath and Ritchey 2012; Cooper et al. 2020; Ritchey and Cooper 2020; Barnett*
724 *et al. 2020; Gilboa and Marlatte 2017; van Kesteren et al. 2012; Preston and Eichenbaum 2013*).
725 We have opted to start with the simplified episodic memory system described in this paper both
726 for reasons of scientific parsimony and also for practical reasons – adding additional neurobiolog-
727 ical details would make the model run too slowly (the current model takes on the order of hours
728 to run on standard computers; adding more complexity would shift this to days or weeks).

729 Just as our model contains some key simplifications, the environment used in the event pro-
730 cessing task is relatively simple and do not capture the full richness of naturalistic events. Some
731 recent studies have explored event graphs with more realistic structure (e.g., *Elman and McRae,*
732 *2019*). The fact that our model can presently only handle one schema substantially limits the com-
733 plexity of the sequences it can process; adding the ability to handle multiple schemas (as discussed
734 above) will help to address this limitation. Also, natural events unfold over multiple timescales. For
735 example, going to the parking lot is an event that involves finding the key, getting to the elevator,
736 etc., but this can be viewed as part of a higher-level event, such as going to an airport. In our simu-
737 lation, events only have one timescale. In general, introducing additional hierarchical structure to
738 the stimuli would enrich the task demands and lead to interesting modeling challenges. For now,
739 we have avoided more complex task environments for computational tractability reasons, but –
740 as computational resources continue to grow – we hope to be able to investigate richer and more
741 realistic task environments going forward. At the same time, we also plan to use the model to

742 address selective retrieval and encoding effects in list-learning studies (e.g., the aforementioned
743 studies showing that surprise boosts encoding; for reviews, see *Frank and Kafkas 2021; Quent*
744 *et al. 2021b*).

745 Another limitation of the model is that the policies explored here (having to do with when
746 episodic memory snapshots are stored and retrieved) do not encompass the full range of ways
747 in which the use of episodic memory can be optimized. For example, in addition to questions
748 about *when* to encode and retrieve, one can consider optimizations of what is stored in memory
749 and how memory is cued. These kinds of optimizations are evident in mnemonic techniques like
750 the method of loci (*Yates, 1966*), which involve considerable recoding of to-be-learned information
751 (to maximize distinctiveness of stored memories) and also structured cuing strategies (to ensure
752 that these distinctive memory traces can be found after they are stored). We think that the kinds
753 of policies explored in this paper (e.g., retrieving more when uncertain, encoding more at the end
754 of an event) fall more on the “automatic” end of the spectrum, as evidenced by the fact that they
755 require no special training and are deployed even in incidental learning situations (e.g., while peo-
756 ple are watching a movie, without specifically trying to remember it; *Chen et al. 2016; Baldassano*
757 *et al. 2017*). As such, these policies seem very different from more complex and deliberate kinds
758 of mnemonic strategies like method of loci that require special training. However, we think that it
759 is best to view our “simple” policies and more complex strategies as falling on a continuum. While
760 the policies we discuss may be deployed automatically in adults, our simulations show that at
761 least some of these policies (e.g., modulating episodic retrieval based on predictive certainty) can
762 be learned through experience, and indeed these strategies might not (yet) be automatic in young
763 children. Furthermore, in principle, there is nothing stopping a model like ours from learning more
764 elaborate strategies given the right kinds of experience and a rich enough action space. Expanding
765 the space of “memory use policies” for our model and exploring how these can be learned is an
766 important future direction for this work (for a resource-rational approach to memory search, see
767 *Zhang et al. 2021*).

768 Lastly, although we have focused on cognitive modeling in this paper, we think that some of
769 our results have implications for machine learning more broadly. For example, most memory-
770 augmented neural networks used in machine learning encode at each time point (*Graves et al.,*
771 *2014, 2016; Ritter et al., 2018; Pritzel et al., 2017*). Our results provide initial evidence that taking
772 episodic “snapshots” too frequently can actually harm performance. Future work can explore the
773 circumstances under which more selective encoding and retrieval policies might lead to improved
774 performance on machine learning benchmarks. Based on our simulations, we expect that these
775 selective policies will be most useful when there is a substantial risk of recalling lure memories
776 that lead to incorrect predictions, and a substantial cost associated with making these incorrect
777 predictions.

778 **Summary**

779 The modeling work presented here builds on a wide range of research showing that episodic mem-
780 ory is a resource that the brain can flexibly draw upon to solve tasks (see, e.g., *Shohamy and Turk-*
781 *Browne 2013; Palombo et al. 2015, 2019; Bakkour et al. 2019; Biderman et al. 2020*). This view
782 implies that, in addition to studying episodic memory using tasks that probe this system in isola-
783 tion, it is also valuable to study how episodic memory is used in more complex situations, in concert
784 with working memory and semantic memory, to solve tasks and make predictions. To engage with
785 findings of these sort, we have leveraged advances in AI that make it possible for models to learn
786 how to use episodic memory – our simulations provide a way of predicting how episodic memory
787 should be deployed to obtain rewards, as a function of the properties of the learning environment.
788 While our understanding of these more complex situations is still at an early stage, our hope is that
789 this model (and others like it, such as the model by *Ritter et al. 2018*) can spur a virtuous cycle of
790 predictions, experiments, and model revision that will bring us to a richer understanding of how
791 the brain uses episodic memory.

792 Methods

793 Episodic retrieval

794 Episodic retrieval in our model is content-based. The retrieval process returns a weighted average
795 of all episodic memories, where the weight of each memory is equal to its activation; to calculate
796 the activation for each memory, the model executes an evidence accumulation process using a
797 leaky competing accumulator (LCA; *Usher and McClelland 2001*), which has been used in other
798 memory models (e.g., *Sederberg et al. 2008; Polyn et al. 2009*). The evidence for a given episodic
799 memory is the cosine similarity between that memory and the current cortical pattern (the cell
800 state of the LSTM). Hence, memories that are similar to the current cortical pattern will have a
801 larger influence on the pattern that gets reinstated. This conceptualization of episodic memory is
802 similar to an attractor network (*Hopfield, 1982; Rolls, 2010*) – each episodic memory serves as an
803 attractor in the space of LSTM cell states, and retrieval moves the LSTM cell state towards those
804 episodic memories.

805 The evidence accumulation process is governed by the episodic memory gate (EM gate) and
806 the level of competition across memories (Figure 1C), which are stored separately from each other
807 (Figure 1B). The EM gate is controlled by the cortical network (Figure 1A, C). The EM gate, in turn,
808 controls whether episodic retrieval happens – a higher EM gate value increases the activation of
809 all memories, and setting the EM gate value to zero turns off episodic retrieval completely (see
810 *Appendix 4* for discussion of other ways that gating can be configured). The level of competition
811 (i.e., lateral inhibition) adjusts the contrast of activations across all memories; making the level of
812 competition higher or lower interpolates between one-winner-take-all recall versus recalling an
813 average of multiple memories. In all of our simulations, we set the level of competition to be well
814 above zero (0.8, to be exact), given the overwhelming evidence that episodic retrieval is competitive
815 (*Anderson and Reder, 1999; Norman and O'Reilly, 2003; Norman, 2010*).

816 Note that, instead of optimizing the LCA parameters to fit empirical results (e.g., as in the work
817 by *Polyn et al. 2009*), we use a neural network that learns to control the level of the EM gate value.
818 As described below, in the *Model training and testing* section, the model's goal is to maximize reward
819 by making correct predictions and avoiding incorrect predictions; the network learns a policy for
820 setting the EM gate value that maximizes the reward it receives. We made several simplifications
821 to the original LCA – in our model, the LCA 1) has no leak; 2) has no noise; and 3) uses the same
822 EM gate value and competition value for all accumulators.

823 Episodic retrieval - Detail

824 At time t , assume the model has n memories. The model first computes the evidence for all of the
825 memories. The evidence for the i -th memory, m_i , is the cosine similarity between the current LSTM
826 cell state pattern c_t and that episodic memory – which is a previously saved LSTM cell state pattern.
827 We denote the evidence for the i -th memory as x_i :

$$x_i = \text{cosine}(c_t, m_i)$$

828 The x_i , for all i , are the input to the evidence accumulation (LCA) process used in our model; the
829 evidence accumulation process has a timescale τ that is faster than t , such that the accumulation
830 process runs to completion within a single time point of the cortical model. The computation at
831 time τ (for $\tau > 0$) is governed by the following formula:

$$w_\tau^i = \text{relu}\left(\alpha x_i - \beta \sum_{j \neq i} w_\tau^j\right)$$

832 w_τ^i is the activation value for the i -th memory at time τ . The activation for the i -th memory is
833 positively related to its evidence, x_i , and is multiplicatively modulated by α , the EM gate value. The
834 i -th memory also receives inhibition from all of the other memories, where the level of inhibition

835 is modulated by the level of competition, β . Finally, the retrieved item at time t , denoted by μ_t , is a
836 combination of all memories, weighted by their activation:

$$\mu_t = \sum_{i=1}^n w_i m_i$$

837 Model training and testing

838 Model training

839 Before the model is used to simulate any particular experiment, it undergoes a *meta-training* phase
840 that is meant to reflect the experience that a person has prior to the experiment. The goal of this
841 meta-training phase is to let the model learn 1) the structure of the task – how situation features
842 control the transition dynamics across states; and 2) a policy for retrieving episodic memories and
843 for making next-state predictions that maximizes the reward it receives. For every epoch of meta-
844 training, it is trained for all three conditions (recent memory, distant memory, and no memory).

845 The model is trained with reinforcement learning. Specifically, the model is rewarded/penalized
846 if its prediction about the next state is correct/incorrect. The model also has the option of saying
847 “don’t know” (implemented as a dedicated output unit) when it is uncertain about what will happen
848 next; if the model says “don’t know”, the reward is zero. The model is trained with the advantage
849 actor-critic (A2C) objective (Mnih et al., 2016). At time t , the model outputs its prediction about the
850 next state, \hat{s}_{t+1} , and an estimate of the state value, v_t . After every event (i.e., a sequence of states of
851 length T), it takes the following policy gradient step to adjust the connection weights for all layers,
852 denoted by θ :

$$\nabla_{\theta} J(\theta) = \nabla \sum_{t=0}^T \log \pi_{\theta}(\hat{s}_{t+1} | s_t) (r_t - v_t)$$

853 This objective makes rewarded actions (next-state predictions) more likely to occur; the above
854 equation shows how this process is modulated by the level of reward prediction error – measured
855 as the difference between the predicted value, v_t , versus the reward at time t , denoted by r_t . We
856 also used entropy regularization on the network output (Grandvalet and Bengio, 2006; Mnih et al.,
857 2016) to encourage exploration in the early phase of the training process.

858 We used the A2C method (Mnih et al., 2016), as it is simple and has been widely used in cognitive
859 modeling (Ritter et al., 2018; Wang et al., 2018). Notably, there is also evidence that an actor-
860 critic style system is implemented in the cortex and basal ganglia (Takahashi et al., 2008). Since
861 pure reinforcement learning is not data-efficient enough, we used supervised initialization during
862 meta-training to help the model develop useful representations (Misra et al., 2017; Nagabandi
863 et al., 2017). Specifically, the model is first trained for 600 epochs to predict the next state and
864 to minimize the cross-entropy loss between the output and the target. During this supervised
865 pre-training phase, the model is only trained on the recent memory condition and the episodic
866 memory module is turned off, so this supervised pre-training does not influence the network’s
867 retrieval policy. Additionally, the “don’t know” output unit is not trained during the supervised pre-
868 training phase – we did this because we want the model to learn its own policy for saying “don’t
869 know”, rather than having one imposed by us. Next, the model is switched to the advantage actor-
870 critic (A2C) objective (Mnih et al., 2016) and trained for another 400 epochs, allowing all weights to
871 be adjusted. The number of training epochs was picked to ensure the learning curves converge.

872 Stimulus representation

873 At time t , the model observes a situation feature, and then it gets a query about which state will be
874 visited next. Specifically, the input vector at time t has four components (see Figure 5): 1) The
875 observed situation feature (sticking with the example in Figure 2, this could be something like
876 “weather”) is encoded as a T -dimensional one-hot vector. T is the total number of situation fea-
877 tures, which (in most simulations) is the same as the number of time points in the event. The t -th

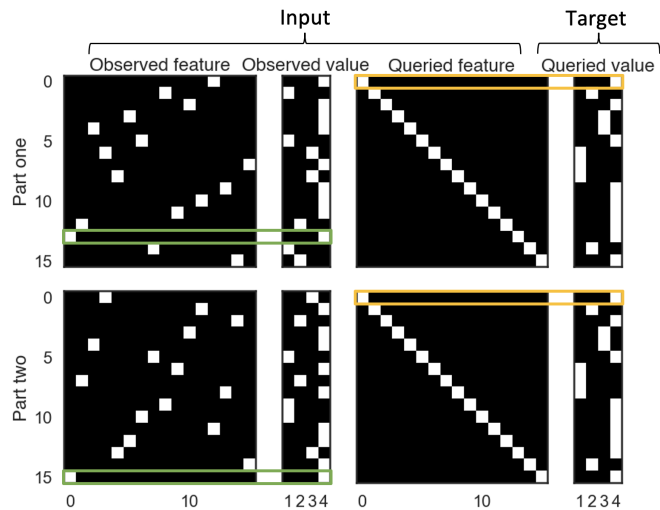


Figure 5. The stimulus representation for the event processing task. In the event processing task, situation features are observed in different, random orders during part 1 and part 2, but queries about those features are presented in the same order during part 1 and part 2. The green boxes in panel A indicate time points where the model observed the value of the first feature (time point 13 during part 1, and time point 15 during part 2). The yellow boxes indicate time points where the model was queried about the value of the first feature (time point zero during both part 1 and part 2).

878 one-hot indicates the situation feature governing the transition at time t . 2) The value of the ob-
879 served situation feature (e.g., learning that the weather is sunny) is encoded as a B -dimensional
880 one-hot vector, where B is the number of possible next states at time t . 3) The queried situation
881 feature is encoded as another T -dimensional one-hot vector (note that querying the model
882 the value of the feature that controls the next-state transition is equivalent to querying the model
883 about the next state, given that there is a 1-to-1 mapping between feature values and states within
884 a time point; see Figure 2A). 4) Finally, the model also receives the current penalty level for incor-
885 rect predictions as a scalar input, which can change across trials. Overall, the input vector at time t
886 is $2T + B + 1$ dimensional. At every time point, there is also a target vector of length B that specifies
887 the value of the queried feature (i.e., the “correct answer” that the model is trying to predict). The
888 model outputs a vector of length $B + 1$: The first B dimensions correspond to specific predictions
889 about the next state, and the last output dimension corresponds to the “don’t know” response.

890 In our simulation, the length of an event is 16 time points (i.e., $T = 16$), and the number of
891 possible states at each time point is 4 (i.e., $B = 4$). Hence the chance level for next-state prediction
892 is $1/4$. Figure 5 illustrates the stimuli provided to the model for a single example trial. Note that
893 the queries (about the next state) are always presented in the same order, so there is a diagonal
894 on the queried feature matrix. This captures the sequential structure of events (e.g., ordering food
895 always happens before eating the food). However, the order in which the situation features are
896 observed is random. As a result, sometimes a feature is queried after it was observed, in which
897 case the model can rely on its working memory to produce the correct prediction, and sometimes
898 a feature is queried before it was observed, in which case the model needs to use episodic memory
899 (if a relevant memory is available) to produce the correct prediction.

900 As discussed above, the input vector specifies the level of penalty (for incorrect prediction) for
901 the current trial. During meta-training, the penalty value was randomly sampled on each trial from
902 the range between 0 and 4. During meta-testing, we evaluated the model using a penalty value
903 of 2 (the average of the penalty values used during training). To understand the effect of penalty
904 on retrieval policy, we also compared the timing of recall in the model when the penalty during
905 meta-testing was low (penalty = 0) vs. high (penalty = 4; Figure 3F).

906 In our simulations, during meta-training, the model only got to observe 70% of the features
907 of the ongoing situation during part 1 of the sequence. This was operationalized by giving each
908 feature a 30% probability of being removed during part 1; for time points where the to-be-observed
909 feature was removed, the model observed a zero vector instead. This “feature removal” during part
910 1 of the sequence made the task more realistic, since – in general – past information does not fully
911 inform what will happen in the future (during meta-testing, we did not remove any observations
912 during part 1; this makes the results graphs easier to interpret, but has no effect on the conclusions
913 reported here).

914 Finally, we wanted to make sure the model could adjust its retrieval time flexibly, instead of
915 learning to always retrieve at a fixed time point (e.g., always retrieve at the third time point). There-
916 fore, during training, we delayed the prediction demand by a random number of time points (from
917 0 to 3). For example, if the amount of delay was 2 in a given trial, then the model observed 2
918 situation features before it received the first query.

919 Model testing

920 During meta-testing (i.e., model evaluation; when simulating a particular experiment), the weights
921 of the cortical part of the model (i.e., all weights pertaining to the LSTM, decision layer, and EM gate)
922 were frozen, but the model was allowed to form new episodic memories. In any given trial (where
923 the model observed several events), new learning of information completely relied on working
924 memory (i.e., model’s recurrent dynamics), episodic memory in the episodic module, and semantic
925 memory encoded in the (frozen) cortical connection weights (instantiating the model’s knowledge
926 of transitions between states and how these transitions are controlled by situation features). The
927 results shown in all of the simulations were obtained by testing the model with new, randomly-
928 generated events, after the initial meta-training phase. While it is theoretically possible that these
929 test events could duplicate events that were encountered during meta-training, exact repeats will
930 be very rare due to the combinatorics of the stimuli (as noted earlier, there are 4^{16} possible se-
931 quences of states within an event). For more information on model parameters, see *Appendix 7*.

932 Decoding the working memory state

933 In Figure 2C, we used a decoding approach to track what information the model was maintaining
934 in working memory over time while it processes an event. This approach allowed us to assess the
935 model’s ability to hold on to observed features after they were observed, and also to detect when
936 features were retrieved from episodic memory and loaded back up into working memory. Our
937 use of decoders here is analogous to the widespread use of multivariate pattern analysis (MVPA)
938 methods to decode the internal states of participants from neuroimaging data (*Haxby et al., 2001*;
939 *Norman et al., 2006*; *Lewis-Peacock and Norman, 2014*) – the only difference is that, here, we
940 applied the decoder to the internal states of the model instead of applying it to brain data.

941 Specifically, we trained classifiers on LSTM cell states during part 1 to decode the feature values
942 over time. Each situation feature was given its own classifier (logistic regression with L2 penalty).
943 For example, if “weather” was one of the situation features, we would train a dedicated “weather”
944 classifier that takes the LSTM cell state and predicts the value of the weather feature for a given
945 time point. To set up the targets for these classifiers for part 1, we labeled all time points before the
946 model observed the feature value as “don’t know”. After a feature value was revealed, we labeled
947 that time point and the following time points with the value of that feature (e.g., if the weather
948 feature value was observed to be “rainy” on time point 4, then time point 4 and all of those that
949 followed until the end of part 1 of the sequence were labeled with the value “rainy”). For part 2 data,
950 we assumed all features were reinstated to the model’s working memory state after the EM gate
951 value peaked. This labeling scheme assumes that 1) observed features are maintained in working
952 memory and 2) episodic recall brings back previously encoded information. These assumptions
953 can be tested by applying the classifier to held-out data. When decoding working memory states
954 during part 1 of the sequence, we used a five-fold cross-validation procedure, and picked the regu-

955 larization parameter with an inner-loop cross-validation. All results were generated using held-out
956 test sets. The average decoding accuracy was 91.58%. Note that, as mentioned above, there is no
957 guarantee that features observed earlier in the sequence will be maintained in the model's work-
958 ing memory. As such, below-ceiling decoding accuracy could reflect either 1) failure to accurately
959 decode the contents of working memory, or 2) the decoder accurately detecting a working memory
960 failure (i.e., that the feature in question has “dropped out” of the model's working memory, despite
961 having been observed earlier in the sequence).

962 Acknowledgments

963 This work was supported by a Multi-University Research Initiative grant awarded to KAN and UH
964 (ONR/DoD N00014-17-1-2961). We are grateful for the feedback we have received from members
965 of the Princeton Computational Memory Lab, the Hasson Lab, and the labs of our MURI collabora-
966 tors Charan Ranganath, Lucia Melloni, Jeffrey Zacks, and Samuel Gershman.

967 Code

968 Github repo: <https://github.com/qihongl/learn-hippo>

969 References

- 970 **Aldous DJ**. Exchangeability and related topics. In: Aldous DJ, Ibragimov IA, Jacod J, editors. *École d'Été de Proba-*
971 *bilités de Saint-Flour XIII — 1983* Springer Berlin Heidelberg; 1985. p. 1–198.
- 972 **Anderson JR**. The adaptive nature of human categorization. *Psychological Review*. 1991; .
- 973 **Anderson JR**, Reder LM. The fan effect: New results and new theories. *Journal of Experimental Psychology:*
974 *General*. 1999; 128(2):186–197.
- 975 **Anderson JR**, Schooler LJ. The adaptive nature of memory. In: Tulving E, editor. *The Oxford Handbook of Memory*
976 Oxford University Press, UK; 2000.p. 557–570.
- 977 **Antony JW**, Hartshorne TH, Pomeroy K, Gureckis TM, Hasson U, McDougale SD, Norman KA. Behavioral, physi-
978 ological, and neural signatures of surprise during naturalistic sports viewing. *Neuron*. 2021 Jan; 109(2):377–
979 390.e7.
- 980 **Axmacher N**, Cohen MX, Fell J, Haupt S, Dümpelmann M, Elger CE, Schlaepfer TE, Lenartz D, Sturm V, Ranganath
981 C. Intracranial EEG correlates of expectancy and memory formation in the human hippocampus and nucleus
982 accumbens. *Neuron*. 2010 Feb; 65(4):541–549.
- 983 **Baddeley A**. The episodic buffer: A new component of working memory? *Trends in Cognitive Sciences*. 2000
984 Nov; 4(11):417–423.
- 985 **Bakkour A**, Palombo DJ, Zylberberg A, Kang YH, Reid A, Verfaellie M, Shadlen MN, Shohamy D. The hippocam-
986 pus supports deliberation during value-based decisions. *eLife*. 2019 Jul; 8.
- 987 **Baldassano C**, Chen J, Zadbood A, Pillow JW, Hasson U, Norman KA. Discovering event structure in continuous
988 narrative perception and memory. *Neuron*. 2017 Aug; 95(3):709–721.e5.
- 989 **Barnett AJ**, Reilly W, Dimsdale-Zucker H, Mizrak E, Reagh Z, Ranganath C. Organization of cortico-hippocampal
990 networks in the human brain. *bioRxiv*. 2020; p. 2020.06.09.142166.
- 991 **Bein O**, Duncan K, Davachi L. Mnemonic prediction errors bias hippocampal states. *Nature Communications*.
992 2020 Jul; 11(1):3451.
- 993 **Ben-Yakov A**, Dudai Y. Constructing realistic engrams: Poststimulus activity of hippocampus and dorsal stria-
994 tum predicts subsequent episodic memory. *The Journal of Neuroscience: the Official Journal of the Society*
995 *for Neuroscience*. 2011 Jun; 31(24):9032–9042.
- 996 **Ben-Yakov A**, Eshel N, Dudai Y. Hippocampal immediate poststimulus activity in the encoding of consecutive
997 naturalistic episodes. *Journal of Experimental Psychology General*. 2013 Nov; 142(4):1255–1263.

- 998 **Ben-Yakov A**, Henson RN. The hippocampal film editor: Sensitivity and specificity to event boundaries in
999 continuous experience. *The Journal of Neuroscience: the Official Journal of the Society for Neuroscience*.
1000 2018 Nov; 38(47):10057–10068.
- 1001 **Biderman N**, Bakkour A, Shohamy D. What are memories for? The hippocampus bridges past experience with
1002 future decisions. *Trends in Cognitive Sciences*. 2020 Jun; 0(0).
- 1003 **Bonasia K**, Sekeres MJ, Gilboa A, Grady CL, Winocur G, Moscovitch M. Prior knowledge modulates the neural
1004 substrates of encoding and retrieving naturalistic events at short and long delays. *Neurobiology of Learning
1005 and Memory*. 2018 Sep; 153(Pt A):26–39.
- 1006 **Botvinick M**, Ritter S, Wang JX, Kurth-Nelson Z, Blundell C, Hassabis D. Reinforcement learning, fast and slow.
1007 *Trends in Cognitive Sciences*. 2019 May; 23(5):408–422.
- 1008 **Brunec IK**, Robin J, Olsen RK, Moscovitch M, Barense MD. Integration and differentiation of hippocampal mem-
1009 ory traces. *Neuroscience and Biobehavioral Reviews*. 2020 Jul; 118:196–208.
- 1010 **Chanales AJH**, Oza A, Favila SE, Kuhl BA. Overlap among spatial memories triggers repulsion of hippocampal
1011 representations. *Current Biology*. 2017 Aug; 27(15):2307–2317.e5.
- 1012 **Chang CHC**, Lazaridi C, Yeshurun Y, Norman KA, Hasson U. Relating the past with the present: Information
1013 integration and segregation during ongoing narrative processing. *Journal of Cognitive Neuroscience*. 2021
1014 May; 33(6):1106–1128.
- 1015 **Chen J**, Cook PA, Wagner AD. Prediction strength modulates responses in human area CA1 to sequence viola-
1016 tions. *Journal of Neurophysiology*. 2015 Aug; 114(2):1227–1238.
- 1017 **Chen J**, Honey CJ, Simony E, Arcaro MJ, Norman KA, Hasson U. Accessing real-life episodic information from
1018 minutes versus hours earlier modulates hippocampal and high-order cortical dynamics. *Cerebral Cortex*.
1019 2016 Aug; 26(8):3428–3441.
- 1020 **Chen PH**, Chen J, Yeshurun Y, Hasson U, Haxby J, Ramadge PJ. A reduced-dimension fMRI shared response
1021 model. *Advances in Neural Information Processing Systems* 28. 2015; .
- 1022 **Chien HYS**, Honey CJ. Constructing and forgetting temporal context in the human cerebral cortex. *Neuron*.
1023 2020 May; 106(4):675–686.e11.
- 1024 **Collins AGE**. Reinforcement learning: Bringing together computation and cognition. *Current Opinion in Behav-
1025 iorral Sciences*. 2019 Oct; 29:63–68.
- 1026 **Cooper RA**, Kurkela KA, Davis SW, Ritchey M. Mapping the organization and dynamics of the posterior medial
1027 network during movie watching. *bioRxiv*. 2020; p. 2020.10.21.348953.
- 1028 **Dauphin YN**, Pascanu R, Gulcehre C, Cho K, Ganguli S, Bengio Y. Identifying and attacking the saddle point
1029 problem in high-dimensional non-convex optimization. *Advances in Neural Information Processing Systems*.
1030 2014; .
- 1031 **Davidow JY**, Foerde K, Galván A, Shohamy D. An upside to reward sensitivity: The hippocampus supports
1032 enhanced reinforcement learning in adolescence. *Neuron*. 2016 Oct; 92(1):93–99.
- 1033 **Dudai Y**. Predicting not to predict too much: How the cellular machinery of memory anticipates the uncertain
1034 future. *Philosophical Transactions of the Royal Society of London Series B, Biological Sciences*. 2009 May;
1035 364(1521):1255–1262.
- 1036 **Dudai Y**, Eisenberg M. Rites of passage of the engram: Reconsolidation and the lingering consolidation hypoth-
1037 esis. *Neuron*. 2004 Sep; 44(1):93–100.
- 1038 **Duncan K**, Sadanand A, Davachi L. Memory's penumbra: Episodic memory decisions induce lingering
1039 mnemonic biases. *Science*. 2012 Jul; 337(6093):485–487.
- 1040 **Duncan K**, Schlichting M. Hippocampal representations as a function of time, subregion, and brain state. *Neu-
1041 robiology of Learning and Memory*. 2018 Sep; 153(Pt A):40–56.
- 1042 **Duncan K**, Semmler A, Shohamy D. Modulating the use of multiple memory systems in value-based decisions
1043 with contextual novelty. *Journal of Cognitive Neuroscience*. 2019 Oct; 31(10):1455–1467.
- 1044 **Duncan K**, Shohamy D. Memory states influence value-based decisions. *Journal of Experimental Psychology:
1045 General*. 2016 Nov; 145(11):1420–1426.

- 1046 **Eichenbaum H.** Time cells in the hippocampus: A new dimension for mapping memories. *Nature Reviews*
1047 *Neuroscience*. 2014 Nov; 15(11):732–744.
- 1048 **Elman JL, McRae K.** A model of event knowledge. *Psychological Review*. 2019 Mar; 126(2):252–291.
- 1049 **Ericsson KA, Kintsch W.** Long-term working memory. *Psychological Review*. 1995 Apr; 102(2):211–245.
- 1050 **Ezzyat Y, Davachi L.** What constitutes an episode in episodic memory? *Psychological Science*. 2011 Feb;
1051 22(2):243–252.
- 1052 **Ezzyat Y, Davachi L.** Neural evidence for representational persistence within events. *The Journal of Neuro-*
1053 *science: the Official Journal of the Society for Neuroscience*. 2021 Jul; .
- 1054 **Favila SE, Chanales AJH, Kuhl BA.** Experience-dependent hippocampal pattern differentiation prevents inter-
1055 ference during subsequent learning. *Nature Communications*. 2016 Apr; 7:11066.
- 1056 **Frank D, Kafkas A.** Expectation-driven novelty effects in episodic memory. *Neurobiology of Learning and*
1057 *Memory*. 2021 May; p. 107466.
- 1058 **Frank D, Kafkas A, Montaldi D.** Experiencing surprise: The temporal dynamics of its impact on memory. *bioRxiv*.
1059 2021 Jul; p. 2020.12.15.422817.
- 1060 **Frank D, Montemurro MA, Montaldi D.** Pattern separation underpins expectation-modulated memory. *The*
1061 *Journal of Neuroscience: the Official Journal of the Society for Neuroscience*. 2020 Apr; 40(17):3455–3464.
- 1062 **Franklin NT, Norman KA, Ranganath C, Zacks JM, Gershman SJ.** Structured Event Memory: A neuro-symbolic
1063 model of event cognition. *Psychological Review*. 2020 Apr; 127(3):327–361.
- 1064 **Gershman SJ.** The adaptive nature of memory. In: Kahana MJ, Wagner AD, editors. *Oxford Handbook of Human*
1065 *Memory*, vol. 700 Oxford University Press, UK; 2021.
- 1066 **Gershman SJ, Blei DM, Niv Y.** Context, learning, and extinction. *Psychological Review*. 2010 Jan; 117(1):197–209.
- 1067 **Gershman SJ, Norman KA, Niv Y.** Discovering latent causes in reinforcement learning. *Current Opinion in*
1068 *Behavioral Sciences*. 2015 Oct; 5:43–50.
- 1069 **Gilboa A, Marlatte H.** Neurobiology of schemas and schema-mediated Memory. *Trends in Cognitive Sciences*.
1070 2017 Aug; 21(8):618–631.
- 1071 **Grandvalet Y, Bengio Y.** Entropy regularization. In: Olivier Chapelle, Bernhard Scholkopf, and Alexander Zien,
1072 editor. *Semi-Supervised Learning*; 2006.
- 1073 **Graves A, Wayne G, Danihelka I.** Neural Turing machines. *arXiv*. 2014 Oct; .
- 1074 **Graves A, Wayne G, Reynolds M, Harley T, Danihelka I, Grabska-Barwińska A, Colmenarejo SG, Grefenstette E,**
1075 **Ramalho T, Agapiou J, Badia AP, Hermann KM, Zwols Y, Ostrovski G, Cain A, King H, Summerfield C, Blunsom**
1076 **P, Kavukcuoglu K, Hassabis D.** Hybrid computing using a neural network with dynamic external memory.
1077 *Nature*. 2016 Oct; 538(7626):471–476.
- 1078 **Greve A, Cooper E, Kaula A, Anderson MC, Henson R.** Does prediction error drive one-shot declarative learning?
1079 *Journal of Memory and Language*. 2017 Jun; 94:149–165.
- 1080 **Greve A, Cooper E, Tibon R, Henson RN.** Knowledge is power: Prior knowledge aids memory for both congru-
1081 ent and incongruent events, but in different ways. *Journal of Experimental Psychology: General*. 2019 Feb;
1082 148(2):325–341.
- 1083 **Griffiths TL, Lieder F, Goodman ND.** Rational use of cognitive resources: Levels of analysis between the com-
1084 putational and the algorithmic. *Topics in Cognitive Science*. 2015 Apr; 7(2):217–229.
- 1085 **Hamilton LS, Huth AG.** The revolution will not be controlled: Natural stimuli in speech neuroscience. *Language,*
1086 *Cognition and Neuroscience*. 2018 Jul; p. 1–10.
- 1087 **Hardt O, Einarsson EÖ, Nader K.** A bridge over troubled water: Reconsolidation as a link between cognitive
1088 and neuroscientific memory research traditions. *Annual Review of Psychology*. 2010 Jan; 61(1):141–167.
- 1089 **Hasselmo ME, Wyble BP.** Free recall and recognition in a network model of the hippocampus: Simulating
1090 effects of scopolamine on human memory function. *Behavioural Brain Research*. 1997 Dec; 89(1-2):1–34.

- 1091 **Hasson U**, Chen J, Honey CJ. Hierarchical process memory: Memory as an integral component of information
1092 processing. *Trends in Cognitive Sciences*. 2015 Jun; 19(6):304–313.
- 1093 **Hasson U**, Nir Y, Levy I, Fuhrmann G, Malach R. Intersubject synchronization of cortical activity during natural
1094 vision. *Science*. 2004 Mar; 303(5664):1634–1640.
- 1095 **Haxby JV**, Gobbini MI, Furey ML, Ishai A, Schouten JL, Pietrini P. Distributed and overlapping representations
1096 of faces and objects in ventral temporal cortex. *Science*. 2001 Sep; 293(5539):2425–2430.
- 1097 **Haxby JV**, Guntupalli JS, Connolly AC, Halchenko YO, Conroy BR, Gobbini MI, Hanke M, Ramadge PJ. A common,
1098 high-dimensional model of the representational space in human ventral temporal cortex. *Neuron*. 2011 Oct;
1099 72(2):404–416.
- 1100 **Haxby JV**, Guntupalli JS, Nastase SA, Feilong M. Hyperalignment: Modeling shared information encoded in
1101 idiosyncratic cortical topographies. *eLife*. 2020 Jun; 9:e56601.
- 1102 **Hochreiter S**, Schmidhuber J. Long short-term memory. *Neural Computation*. 1997 Nov; 9(8):1735–1780.
- 1103 **Holdstock JS**, Mayes AR, Roberts N, Cezayirli E, Isaac CL, O'Reilly RC, Norman KA. Under what conditions is
1104 recognition spared relative to recall after selective hippocampal damage in humans? *Hippocampus*. 2002;
1105 12(3):341–351.
- 1106 **Hopfield JJ**. Neural networks and physical systems with emergent collective computational abilities. *Proceed-*
1107 *ings of the National Academy of Sciences of the United States of America*. 1982 Apr; 79(8):2554–2558.
- 1108 **Hulbert JC**, Norman KA. Neural differentiation tracks improved recall of competing memories following inter-
1109 leaved study and retrieval practice. *Cerebral Cortex*. 2015 Oct; 25(10):3994–4008.
- 1110 **Kafkas A**, Montaldi D. Expectation affects learning and modulates memory experience at retrieval. *Cognition*.
1111 2018 Nov; 180:123–134.
- 1112 **Kafkas A**, Montaldi D. Striatal and midbrain connectivity with the hippocampus selectively boosts memory for
1113 contextual novelty. *Hippocampus*. 2015 Nov; 25(11):1262–1273.
- 1114 **Kahana MJ**. *Foundations of Human Memory*. Oxford University Press, USA; 2012.
- 1115 **van Kesteren MTR**, Ruiters DJ, Fernández G, Henson RN. How schema and novelty augment memory formation.
1116 *Trends in Neurosciences*. 2012 Apr; 35(4):211–219.
- 1117 **Ketz N**, Morkonda SG, O'Reilly RC. Theta coordinated error-driven learning in the hippocampus. *PLOS Compu-*
1118 *tational Biology*. 2013 Jun; 9(6):e1003067.
- 1119 **Kim G**, Norman KA, Turk-Browne NB. Neural differentiation of incorrectly predicted memories. *The Journal of*
1120 *Neuroscience: the Official Journal of the Society for Neuroscience*. 2017 Feb; 37(8):2022–2031.
- 1121 **Kingma DP**, Ba J. Adam: A method for stochastic optimization. *arXiv*. 2014 Dec; .
- 1122 **Koster R**, Chadwick MJ, Chen Y, Berron D, Banino A, Düzel E, Hassabis D, Kumaran D. Big-loop recurrence
1123 within the hippocampal system supports integration of information across episodes. *Neuron*. 2018 Sep;
1124 99(6):1342–1354.e6.
- 1125 **Kumar M**, Ellis CT, Lu Q, Zhang H, Capotă M, Willke TL, Ramadge PJ, Turk-Browne NB, Norman KA. BrainIAK
1126 tutorials: User-friendly learning materials for advanced fMRI analysis. *PLoS computational biology*. 2020 Jan;
1127 16(1):e1007549.
- 1128 **Kumar M**, Michael Anderson, Antony J, Baldassano C, Brooks PP, Cai MB, Chen PHC, Ellis CT, Henselman-
1129 Petrusek G, Huberdeau D, Hutchinson JB, Li PY, Lu Q, Manning JR, Mennen AC, Nastase SA, Richard H, Schapiro
1130 AC, Schuck NW, Shvartsman M, et al. BrainIAK: The brain imaging analysis kit. *OSF Preprints*. 2020; .
- 1131 **Kumaran D**, Maguire EA. An unexpected sequence of events: Mismatch detection in the human hippocampus.
1132 *PLoS Biology*. 2006 Nov; 4(12):e424.
- 1133 **Kumaran D**, Maguire EA. Match–mismatch processes underlie human hippocampal responses to associa-
1134 tive novelty. *The Journal of Neuroscience: the Official Journal of the Society for Neuroscience*. 2007 Aug;
1135 27(32):8517–8524.
- 1136 **Lewis-Peacock JA**, Norman KA. Competition between items in working memory leads to forgetting. *Nature*
1137 *Communications*. 2014 Dec; 5:5768.

- 1138 Li Y, Yosinski J, Clune J, Lipson H, Hopcroft J. Convergent learning: Do different neural networks learn the same
1139 representations? *Proceedings of Machine Learning Research*. 2015; 44:196–212.
- 1140 Lieder F, Griffiths TL. Resource-rational analysis: Understanding human cognition as the optimal use of limited
1141 computational resources. *The Behavioral and Brain Sciences*. 2019 Feb; 43:e1.
- 1142 Long NM, Lee H, Kuhl BA. Hippocampal mismatch signals are modulated by the strength of neural predic-
1143 tions and their similarity to outcomes. *The Journal of Neuroscience: the Official Journal of the Society for*
1144 *Neuroscience*. 2016 Dec; 36(50):12677–12687.
- 1145 Lu Q, Chen PH, Pillow JW, Ramadge PJ, Norman KA, Hasson U. Shared representational geometry across neural
1146 networks. *arXiv*. 2018 Nov; .
- 1147 McClelland JL. Incorporating rapid neocortical learning of new schema-consistent information into comple-
1148 mentary learning systems theory. *Journal of Experimental Psychology: General*. 2013 Nov; 142(4):1190–
1149 1210.
- 1150 McClelland JL, McNaughton BL, Lampinen AK. Integration of new information in memory: new insights from
1151 a complementary learning systems perspective. *Philosophical Transactions of the Royal Society B*. 2020;
1152 375(1799):20190637.
- 1153 McClelland JL, McNaughton BL, O'Reilly RC. Why there are complementary learning systems in the hippocam-
1154 pus and neocortex: Insights from the successes and failures of connectionist models of learning and memory.
1155 *Psychological Review*. 1995 Jul; 102(3):419–457.
- 1156 McClelland JL, Rogers TT. The parallel distributed processing approach to semantic cognition. *Nature Reviews*
1157 *Neuroscience*. 2003 Apr; 4(4):310–322.
- 1158 Meng Q, Chen W, Zheng S, Ye Q, Liu TY. Optimizing neural networks in the equivalent class space. *arXiv*. 2018
1159 Feb; .
- 1160 Michelmann S, Price AR, Aubrey B, Strauss CK, Doyle WK, Friedman D, Dugan PC, Devinsky O, Devore S,
1161 Flinker A, Hasson U, Norman KA. Moment-by-moment tracking of naturalistic learning and its underlying
1162 hippocampo-cortical interactions. *Nature Communications*. 2021 Sep; 12(1):1–15.
- 1163 Misra D, Langford J, Artzi Y. Mapping instructions and visual observations to actions with reinforcement learn-
1164 ing. *arXiv*. 2017 Apr; .
- 1165 Mnih V, Badia AP, Mirza M, Graves A, Lillicrap TP, Harley T, Silver D, Kavukcuoglu K. Asynchronous methods for
1166 deep reinforcement learning. *arXiv*. 2016 Feb; .
- 1167 Nagabandi A, Kahn G, Fearing RS, Levine S. Neural network dynamics for model-based deep reinforcement
1168 learning with model-free fine-tuning. *arXiv*. 2017 Aug; .
- 1169 Nastase SA, Gazzola V, Hasson U, Keysers C. Measuring shared responses across subjects using intersubject
1170 correlation. *Social Cognitive and Affective Neuroscience*. 2019 Aug; 14(6):667–685.
- 1171 Nastase SA, Goldstein A, Hasson U. Keep it real: Rethinking the primacy of experimental control in cognitive
1172 neuroscience. *NeuroImage*. 2020 Aug; 222:117254.
- 1173 Norman KA. How hippocampus and cortex contribute to recognition memory: revisiting the complementary
1174 learning systems model. *Hippocampus*. 2010 Nov; 20(11):1217–1227.
- 1175 Norman KA, Detre G, Polyn SM. Computational Models of Episodic Memory. In: Sun R, editor. *The Cambridge*
1176 *Handbook of Computational Psychology* Cambridge Handbooks in Psychology, Cambridge University Press;
1177 2008.p. 189–225.
- 1178 Norman KA, O'Reilly RC. Modeling hippocampal and neocortical contributions to recognition memory: A
1179 complementary-learning-systems approach. *Psychological Review*. 2003 Oct; 110(4):611–646.
- 1180 Norman KA, Polyn SM, Detre GJ, Haxby JV. Beyond mind-reading: Multi-voxel pattern analysis of fMRI data.
1181 *Trends in Cognitive Sciences*. 2006 Sep; 10(9):424–430.
- 1182 Palombo DJ, Hayes SM, Reid AG, Verfaellie M. Hippocampal contributions to value-based learning: Converging
1183 evidence from fMRI and amnesia. *Cognitive, Affective & Behavioral Neuroscience*. 2019 Jun; 19(3):523–536.
- 1184 Palombo DJ, Keane MM, Verfaellie M. How does the hippocampus shape decisions? *Neurobiology of Learning*
1185 *and Memory*. 2015 Nov; 125:93–97.

- 1186 **Paszke A**, Gross S, Chintala S, Chanan G, Yang E, DeVito Z, Lin Z, Desmaison A, Antiga L, Lerer A. Automatic
1187 differentiation in PyTorch; 2017.
- 1188 **Paszke A**, Gross S, Massa F, Lerer A, Bradbury J, Chanan G, Killeen T, Lin Z, Gimelshein N, Antiga L, Desmaison
1189 A, Köpf A, Yang E, DeVito Z, Raison M, Tejani A, Chilamkurthy S, Steiner B, Fang L, Bai J, et al. PyTorch: An
1190 imperative style, high-performance deep learning library. arXiv. 2019 Dec; .
- 1191 **Patil A**, Duncan K. Lingering cognitive states shape fundamental mnemonic abilities. *Psychological Science*.
1192 2018 Jan; 29(1):45–55.
- 1193 **Pine A**, Sadeh N, Ben-Yakov A, Dudai Y, Mendelsohn A. Knowledge acquisition is governed by striatal prediction
1194 errors. *Nature Communications*. 2018 Apr; 9(1):1673.
- 1195 **Pitman J**. *Combinatorial Stochastic Processes: Ecole d'Été de Probabilités de Saint-Flour XXXII – 2002*. Picard J,
1196 editor, Springer, Berlin, Heidelberg; 2006.
- 1197 **Polyn SM**, Norman KA, Kahana MJ. A context maintenance and retrieval model of organizational processes in
1198 free recall. *Psychological Review*. 2009 Jan; 116(1):129–156.
- 1199 **Preston AR**, Eichenbaum H. Interplay of hippocampus and prefrontal cortex in memory. *Current Biology*. 2013
1200 Sep; 23(17):R764–73.
- 1201 **Pritzel A**, Uria B, Srinivasan S, Badia AP, Vinyals O, Hassabis D, Wierstra D, Blundell C. Neural episodic control.
1202 *Proceedings of Machine Learning Research*. 2017; 70:2827–2836.
- 1203 **Quent JA**, Greve A, Henson R. Shape of U: The relationship between object-location memory and expectedness.
1204 PsyArXiv. 2021 May; .
- 1205 **Quent JA**, Henson RN, Greve A. A predictive account of how novelty influences declarative memory. *Neurobi-
1206 ology of Learning and Memory*. 2021 Mar; 179:107382.
- 1207 **Radvansky GA**, Krawietz SA, Tamplin AK. Walking through doorways causes forgetting: Further explorations.
1208 *Quarterly Journal of Experimental Psychology*. 2011 Aug; 64(8):1632–1645.
- 1209 **Ranganath C**, Ritchey M. Two cortical systems for memory-guided behaviour. *Nature Reviews Neuroscience*.
1210 2012 Oct; 13(10):713–726.
- 1211 **Raposo D**, Ritter S, Santoro A, Wayne G, Weber T, Botvinick M, van Hasselt H, Song F. Synthetic returns for
1212 long-term credit assignment. arXiv. 2021 Feb; .
- 1213 **Reagh ZM**, Delarazan AI, Garber A, Ranganath C. Aging alters neural activity at event boundaries in the hip-
1214 pocampus and Posterior Medial network. *Nature Communications*. 2020 Aug; 11(1):3980.
- 1215 **Richmond LL**, Zacks JM. Constructing experience: Event models from perception to action. *Trends in Cognitive
1216 Sciences*. 2017 Dec; 21(12):962–980.
- 1217 **Ritchey M**, Cooper RA. Deconstructing the posterior medial episodic network. *Trends in Cognitive Sciences*.
1218 2020 Jun; 24(6):451–465.
- 1219 **Ritter S**. *Meta-reinforcement Learning with Episodic Recall: An Integrative Theory of Reward-Driven Learning*.
1220 PhD thesis, Princeton University; 2019.
- 1221 **Ritter S**, Wang JX, Kurth-Nelson Z, Jayakumar SM, Blundell C, Pascanu R, Botvinick M. Been there, done that:
1222 Meta-Learning with episodic recall. *Proceedings of the International Conference on Machine Learning*. 2018;
1223 .
- 1224 **Ritvo VJH**, Turk-Browne NB, Norman KA. Nonmonotonic plasticity: How memory retrieval drives learning.
1225 *Trends in Cognitive Sciences*. 2019 Sep; 23(9):726–742.
- 1226 **Rogers TT**, McClelland JL. *Semantic cognition: A parallel distributed processing approach*, vol. 425. Cambridge,
1227 MA, US: MIT Press Semantic cognition; 2004.
- 1228 **Rolls ET**. *Attractor networks*. *Wiley Interdisciplinary Reviews: Cognitive Science*. 2010 Jan; 1(1):119–134.
- 1229 **Rouhani N**, Norman KA, Niv Y. Dissociable effects of surprising rewards on learning and memory. *Journal of
1230 Experimental Psychology: Learning, Memory, and Cognition*. 2018 Sep; 44(9):1430–1443.

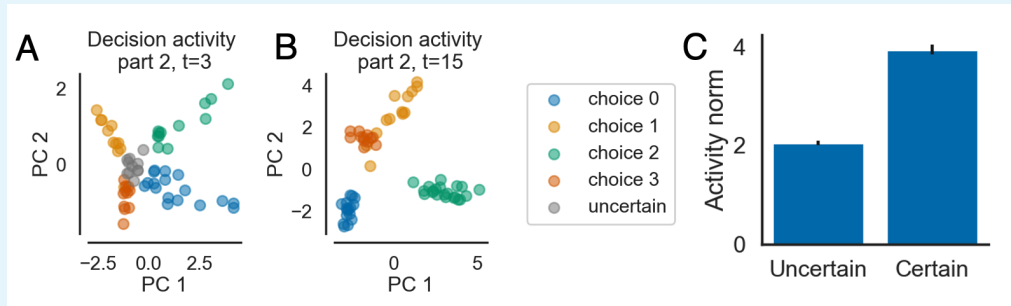
- 1231 **Rouhani N**, Norman KA, Niv Y, Bornstein AM. Reward prediction errors create event boundaries in memory.
1232 *Cognition*. 2020 Oct; 203:104269.
- 1233 **Rumelhart DE**, McClelland JL, Group PR, Others. Parallel distributed processing, vol. 1. MIT press Cambridge,
1234 MA; 1987.
- 1235 **Saxe AM**, McClelland JL, Ganguli S. Exact solutions to the nonlinear dynamics of learning in deep linear neural
1236 networks. *International Conference on Learning Representations*. 2014; .
- 1237 **Saxe AM**, McClelland JL, Ganguli S. A mathematical theory of semantic development in deep neural networks.
1238 *Proceedings of the National Academy of Sciences of the United States of America*. 2019 Jun; 116(23):11537–
1239 11546.
- 1240 **Schapiro AC**, Kustner LV, Turk-Browne NB. Shaping of object representations in the human medial temporal
1241 lobe based on temporal regularities. *Current Biology*. 2012 Sep; 22(17):1622–1627.
- 1242 **Schapiro AC**, Rogers TT, Cordova NI, Turk-Browne NB, Botvinick MM. Neural representations of events arise
1243 from temporal community structure. *Nature Neuroscience*. 2013 Apr; 16(4):486–492.
- 1244 **Schapiro AC**, Turk-Browne NB, Botvinick MM, Norman KA. Complementary learning systems within the hip-
1245 pocampus: A neural network modelling approach to reconciling episodic memory with statistical learn-
1246 ing. *Philosophical Transactions of the Royal Society of London Series B, Biological sciences*. 2017 Jan;
1247 372(1711):20160049.
- 1248 **Schapiro AC**, Turk-Browne NB, Norman KA, Botvinick MM. Statistical learning of temporal community structure
1249 in the hippocampus. *Hippocampus*. 2016 Jan; 26(1):3–8.
- 1250 **Schlichting ML**, Mumford JA, Preston AR. Learning-related representational changes reveal dissociable inte-
1251 gration and separation signatures in the hippocampus and prefrontal cortex. *Nature Communications*. 2015
1252 Aug; 6:8151.
- 1253 **Sederberg PB**, Howard MW, Kahana MJ. A context-based theory of recency and contiguity in free recall. *Psy-
1254 chological Review*. 2008 Oct; 115(4):893–912.
- 1255 **Sherman BE**, Turk-Browne NB. Statistical prediction of the future impairs episodic encoding of the present.
1256 *Proceedings of the National Academy of Sciences of the United States of America*. 2020 Sep; 117(37):22760–
1257 22770.
- 1258 **Shohamy D**, Turk-Browne NB. Mechanisms for widespread hippocampal involvement in cognition. *Journal of
1259 Experimental Psychology: General*. 2013 Nov; 142(4):1159–1170.
- 1260 **Simony E**, Honey CJ, Chen J, Lositsky O, Yeshurun Y, Wiesel A, Hasson U. Dynamic reconfiguration of the default
1261 mode network during narrative comprehension. *Nature Communications*. 2016 Jul; 7:12141.
- 1262 **Sonkusare S**, Breakspear M, Guo C. Naturalistic stimuli in neuroscience: Critically acclaimed. *Trends in Cogni-
1263 tive Sciences*. 2019 Aug; 23(8):699–714.
- 1264 **Stachenfeld KL**, Botvinick MM, Gershman SJ. The hippocampus as a predictive map. *Nature Neuroscience*.
1265 2017 Nov; 20(11):1643–1653.
- 1266 **Stawarczyk D**, Bezdek MA, Zacks JM. Event representations and predictive processing: The role of the midline
1267 default network core. *Topics in Cognitive Science*. 2019 Sep; 30:1345.
- 1268 **van Strien NM**, Cappaert NLM, Witter MP. The anatomy of memory: An interactive overview of the
1269 parahippocampal-hippocampal network. *Nature Reviews Neuroscience*. 2009 Apr; 10(4):272–282.
- 1270 **Sutton RS**, Barto AG. Reinforcement learning: An introduction. MIT press; 2018.
- 1271 **Takahashi Y**, Schoenbaum G, Niv Y. Silencing the critics: Understanding the effects of cocaine sensitization on
1272 dorsolateral and ventral striatum in the context of an actor/critic model. *Frontiers in Neuroscience*. 2008 Jul;
1273 2(1):86–99.
- 1274 **Usher M**, McClelland JL. The time course of perceptual choice: The leaky, competing accumulator model. *Psy-
1275 chological Review*. 2001 Jul; 108(3):550–592.
- 1276 **Wang JX**, Kurth-Nelson Z, Kumaran D, Tirumala D, Soyer H, Leibo JZ, Hassabis D, Botvinick M. Prefrontal cortex
1277 as a meta-reinforcement learning system. *Nature Neuroscience*. 2018 Jun; 21(6):860–868.

- 1278 **Wang SH**, Morris RGM. Hippocampal-neocortical interactions in memory formation, consolidation, and recon-
1279 solidation. *Annual Review of Psychology*. 2010; 61:49–79, C1–4.
- 1280 **Wayne G**, Hung CC, Amos D, Mirza M, Ahuja A, Grabska-Barwinska A, Rae J, Mirowski P, Leibo JZ, Santoro A,
1281 Gemici M, Reynolds M, Harley T, Abramson J, Mohamed S, Rezende D, Saxton D, Cain A, Hillier C, Silver D,
1282 et al. Unsupervised predictive memory in a goal-directed agent. *arXiv*. 2018 Mar; .
- 1283 **Whittington JCR**, Muller TH, Mark S, Chen G, Barry C, Burgess N, Behrens TEJ. The Tolman-Eichenbaum Ma-
1284 chine: Unifying space and relational memory through generalization in the hippocampal formation. *Cell*.
1285 2020 Nov; 183(5):1249–1263.e23.
- 1286 **Yates FA**. *The art of memory*. Chicago: University of Chicago Press; 1966.
- 1287 **Yonelinas AP**. The Nature of recollection and familiarity: A review of 30 years of research. *Journal of Memory*
1288 *and Language*. 2002 Apr; 46(3):441–517.
- 1289 **Zacks JM**. Event perception and memory. *Annual Review of Psychology*. 2020 Jan; 71:165–191.
- 1290 **Zacks JM**, Kurby CA, Eisenberg ML, Haroutunian N. Prediction error associated with the perceptual segmenta-
1291 tion of naturalistic events. *Journal of Cognitive Neuroscience*. 2011 Dec; 23(12):4057–4066.
- 1292 **Zacks JM**, Speer NK, Swallow KM, Braver TS, Reynolds JR. Event perception: A mind-brain perspective. *Psycho-*
1293 *logical Bulletin*. 2007 Mar; 133(2):273–293.
- 1294 **Zhang Q**, Griffiths T, Norman K. Optimal policies for free recall. *PsyArXiv*. 2021 Apr; .

1295 **Appendix 1**

1296

The internal representation of the decision layer



1297

1298

1299

1300

1301

1302

1303

1304

1305

1306

1307

1308

1309

Appendix 1 Figure 1. How certainty is represented in the model's activity patterns. Panels A and B show the neural activity patterns from the decision layer in the distant memory (DM) condition, projected onto the first two principal components. Each point corresponds to the pattern of neural activity for a trial at a particular time point. We colored the points based on the output (i.e., "choice") of the model, which represents the model's belief about which state will happen next. Patterns that subsequently led to "don't know" responses are colored in grey. Panel A shows an early time point with substantial uncertainty (a large number of "don't know" responses). Panel B shows the last time point of this event, where the model has lower uncertainty. Panel C shows the average L2 norm of states that led to "don't know" responses (uncertain) versus states that led to specific next-state predictions (certain); the errorbars indicate 1SE across 15 models. States corresponding to "don't know" responses are clustered in the center of the activation space, with a lower L2 norm.

1310

1311

1312

1313

1314

1315

1316

1317

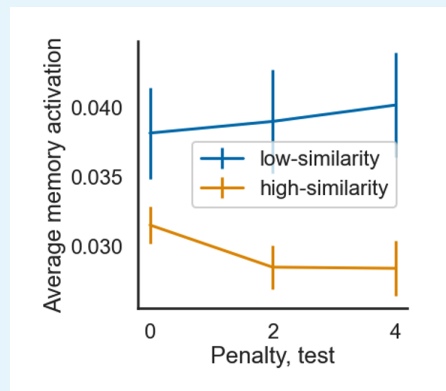
To explore how neural activity patterns in the decision layer differed as a function of certainty, we plotted the activity patterns as a function of the action taken by the model (i.e., whether it predicted one of the four upcoming states, or whether it used the "don't know" response). Figure 1 shows the results of this analysis: Uncertain states are approximately clustered near the center of the activation space (with a lower L2 norm) while other responses are farther away, which indicates that uncertainty in our model is represented by the absence of evidence towards any particular choice. Importantly, this difference in activity patterns is not built-in to the model – it simply emerges during training.

1318 Appendix 2

1319 Effects of event similarity on retrieval policy

1320 In this simulation, we studied how the similarity of event memories in the training environ-
1321 nment affects retrieval policy. To manipulate memory similarity, we varied the proportion of
1322 shared situation feature values across events during training. In the low-similarity condition,
1323 the similarity between the distractor situation (i.e., situation A; see Figure 2 in the main text)
1324 and the target situation was constrained to be less than 40%, so target memories and lures
1325 were relatively easy to distinguish. In the high-similarity condition, the similarity between
1326 the distractor situation and the target situation was constrained to fall between 35% and
1327 90%. We used a rejection sampling approach to implement these similarity bounds – during
1328 stimulus generation, we kept generating distractor situations until they fell within the sim-
1329 ilarity bounds with respect to the target sequence. Otherwise, the simulation parameters
1330 were the same as the parameters that were used in the main text.

1331 In the high-similarity condition, target and lure memories were more confusable, and
1332 thus the risk of lure recall was higher. In light of this, we expected that the model would
1333 adopt a more conservative retrieval policy (i.e., retrieving less) in the high-similarity condi-
1334 tion. We also expected that this effect would be stronger when the penalty is high; when
1335 the penalty is low, there is less of a cost for recalling the lure memory, and thus less of a
1336 reason to refrain from episodic retrieval in the high-similarity condition.



1337 **Appendix 2 Figure 1.** Memory activation during part 2 (averaged over time) in the DM condition, for
1338 models trained in low vs. high event-similarity environments and tested with penalty values that were
1339 low (penalty = 0), moderate (penalty = 2), or high (penalty = 4). The model recalls less when similarity
1340 is high (vs. low), and this effect is larger for higher penalty values. The errorbars indicate 1SE across
1341 15 models.
1342

1344 We compared the model's behavior as a function of penalty and similarity. For the
1345 penalty manipulation, each model was trained on a range of penalty values from 0 to 4,
1346 then tested on low (0), moderate (2), and high (4) penalty values. Figure 1 shows the aver-
1347 age level of memory activation in each of the conditions. As expected, memory activation is
1348 lower in the high-similarity condition, especially when the penalty is high. Notably, increas-
1349 ing penalty reduces memory activation in the high-similarity condition (where the risk of
1350 false recall is high) but it does not have this effect in the low-similarity condition (where the
1351 risk of false recall is low).

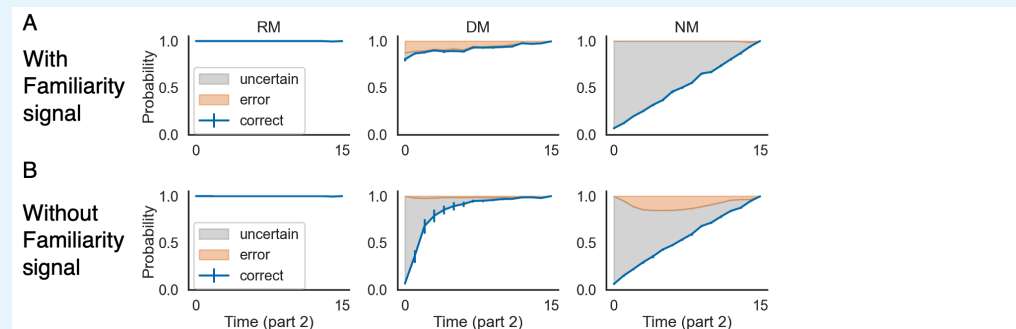
1352 Appendix 3

1353 Effects of familiarity on retrieval policy

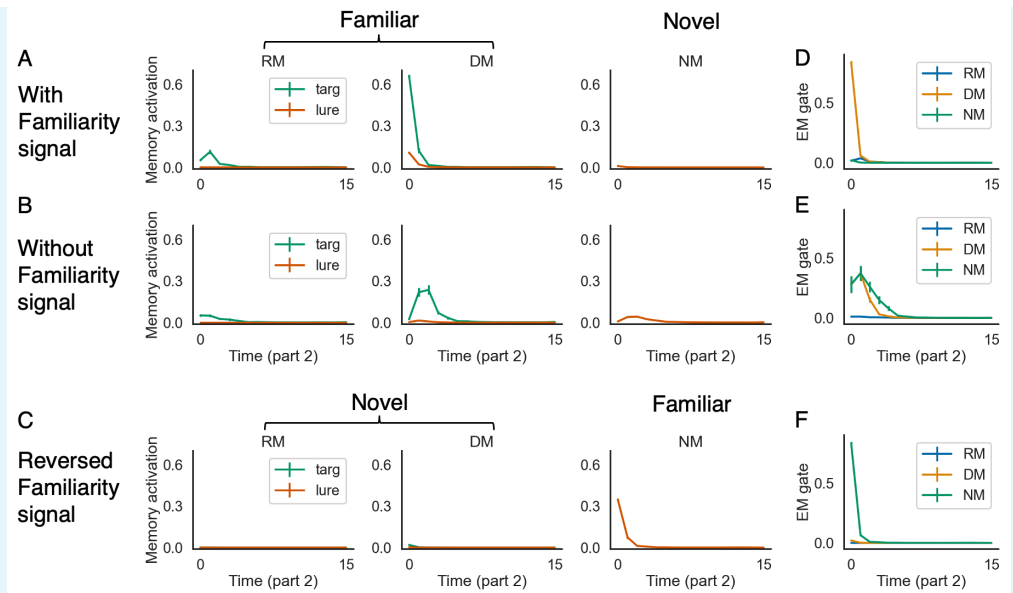
1354 Prior work has demonstrated that cortex is capable of computing a familiarity signal on its
1355 own (i.e., without hippocampus) that discriminates between previously-encountered and
1356 novel stimuli (Yonelinas, 2002; Norman and O'Reilly, 2003; Norman, 2010; Holdstock et al.,
1357 2002). In this section, we study how this familiarity signal can support episodic retrieval pol-
1358 icy. Relevant to this point, several recent studies have found that encountering a familiar
1359 stimulus can temporarily shift the hippocampus into a “retrieval mode” where it is more
1360 likely to retrieve episodic memories in response to available retrieval cues (Duncan et al.,
1361 2012; Duncan and Shohamy, 2016; Duncan et al., 2019; Patil and Duncan, 2018; Hasselmo
1362 and Wyble, 1997). Here, we assess whether our model can provide a resource-rational ac-
1363 count of these “retrieval mode” findings.

1364 Intuitively, familiarity can guide episodic retrieval policy by providing an indication of
1365 whether a relevant episodic memory is available. If an item is unfamiliar, this signals that it
1366 is unlikely that relevant episodic memories exist, hence the expected benefit of retrieving
1367 from episodic memory is low (if there are no relevant episodic memories, episodic retrieval
1368 can only yield irrelevant memories, which lead to incorrect predictions); and if an item is
1369 familiar, this signals that relevant episodic memories are likely to exist and hence the ben-
1370 efits of retrieving from episodic memory are higher. These points suggest that the model
1371 would benefit from a policy whereby it adopts a more liberal criterion for consulting episodic
1372 memory when stimuli are familiar as opposed to novel.

1373 To test this, we ran simulations where we presented a “ground truth” familiarity signal
1374 to the model during part 2 of the sequence. The familiarity signal was presented using an
1375 additional, dedicated input unit (akin to how we present penalty information to the model).
1376 Specifically, during part 2, if the ongoing situation had been observed before (as was the
1377 case in the RM and DM conditions), the familiarity signal was set to one. In contrast, if the
1378 ongoing situation was novel (as was the case in the NM condition), then the familiarity
1379 signal was set to negative one. Before part 2, the familiarity signal was set to zero (an uninfor-
1380 mative value). Other than these changes, the parameters of this simulation were the same
1381 as the other simulations. The model was tested on penalty value of 2 – the average of the
1382 training range. Note that our treatment of the familiarity signal here deliberately glosses
1383 over the question of how this signal is generated, as this question is addressed in detail in
1384 other models (e.g., Norman and O'Reilly 2003); our intent here is to understand the con-
1385 sequences of having a familiarity signal (however it might be generated) for the model's
1386 episodic retrieval policy.



1387 **Appendix 3 Figure 1. The familiarity signal can improve prediction.** Next-state prediction
1388 performance for models with (A) vs. without (B) access to the familiarity signal. With the familiarity
1389 signal (A), the model shows 1) higher levels of correct prediction in the DM condition, and 2) a reduced
1390 error rate in the NM condition. The errorbars indicate 1SE across 15 models.
1392



1393

1394

1395

1396

1397

1398

1399

1400

1401

1402

1403

1404

1405

1406

1407

1408

1409

1410

1411

1412

1413

1414

1415

1416

1417

1418

1419

1420

1421

1422

1423

1424

1425

1426

1427

Appendix 3 Figure 2. Episodic retrieval is modulated by familiarity. This figure shows the memory activation and EM gate values over time for three conditions: 1) with the familiarity signal (A, D), 2) without the familiarity signal (B, E), and 3) with a reversed (opposite) familiarity signal at test (C, F). With the familiarity signal (A), the model shows higher levels of recall in the DM condition, and suppresses recall even further in the NM condition, compared to the model without the familiarity signal (B). This is due to the influence of the EM gate – the model with the familiarity signal retrieves immediately in the DM condition, and turns off episodic retrieval almost completely in the NM condition (D). Note also that levels of episodic retrieval in the RM condition stay low, even with the familiarity signal (see text for discussion). Finally, parts C and F show that reversing the familiarity signal at test suppresses recall in the DM condition and boosts recall in the NM condition. The errorbars indicate 1SE across 15 models.

Figures 1 and 2 illustrate prediction performance, memory activation, and EM gate values for models with and without the familiarity signal. When the model has access to a veridical familiarity signal (+1 for RM and DM, -1 for NM), it opens the EM gate immediately and strongly in the DM condition (Figure 2D - DM), leading to higher activation of both the target memory and the lure (Figure 2A - DM) in the DM condition, relative to models without the familiarity signal (Figure 2B - DM). Behaviorally, models with the familiarity signal show both a higher correct prediction rate and a slightly higher error rate in the DM condition, compared to models without the familiarity signal (Figure 1A vs. B - DM). This slight increase in errors occurs because, when the model retrieves immediately from episodic memory during part 2, the model (in some cases) has not yet made enough observations to distinguish the target and the lure. In the NM condition, with the familiarity signal, the model keeps the EM gate almost completely shut (Figure 2D - NM). Consequently, the level of memory activation stays very low in the NM condition (Figure 2A - NM), which reduces the error rate in the NM condition to zero (Figure 1A - NM). The RM condition is an interesting case: Previously (see Figure 3 in the main text), we found that the model refrained from episodic memory retrieval in the RM condition; we found that the same pattern is present here, even when we make a familiarity signal available to the model: EM gate and memory activation values are both very low (Figure 2A, D - RM), similar to models without access to the familiarity signal (Figure 2B, E - RM). This shows that model does not always retrieve from episodic memory when given a high familiarity signal – in this case, the presence of relevant information in working memory (which suppresses episodic retrieval) “overrides” the presence of the familiarity signal (which enhances episodic retrieval in the DM condition).

Finally, we can trick the model into reversing its retrieval policy by reversing the famil-

1428
1429
1430
1431
1432
1433
1434
1435
1436
1437
1438
1439

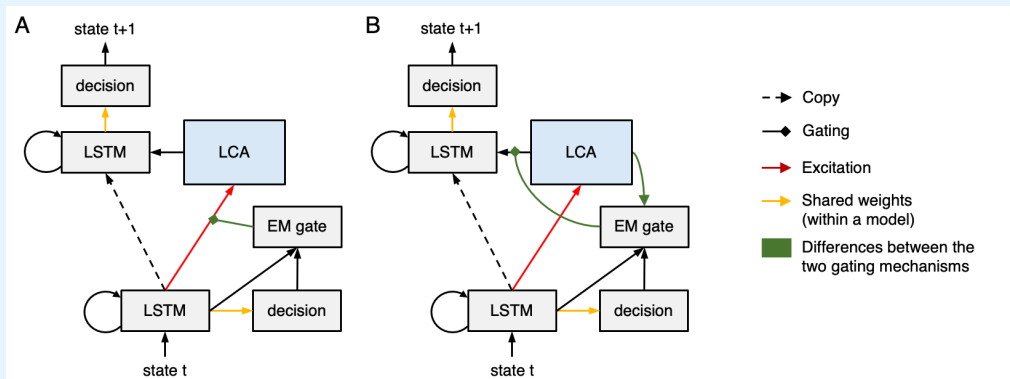
ilarity signal at test (Figure 2C, F). In this condition, the (reversed) signal indicates that the ongoing situation is novel (-1) in the RM and the DM condition, and the ongoing situation is familiar (+1) in the NM condition. As a result, the model suppresses episodic retrieval in the RM and DM conditions, and recalls lures in the NM condition.

Overall, the results of this simulation show that our model is able to use a familiarity signal to inform its retrieval policy in the service of predicting upcoming states. Consistent with empirical results (*Duncan et al., 2012; Duncan and Shohamy, 2016; Duncan et al., 2019; Patil and Duncan, 2018; Hasselmo and Wyble, 1997*), we found that the model retrieves more from episodic memory when the ongoing situation is familiar, unless the model has low uncertainty about the upcoming state. These modeling results provide a resource-rational account of why familiarity leads to enhanced episodic retrieval.

1440 Appendix 4

1441 Alternative configurations of episodic memory gating

1442 In the simulations described in the main text, the EM gate controls the input into the EM
1443 system. An alternative way of accomplishing gating is to place the gate *after* the EM module
1444 (LCA), so it controls the flow of activation from the EM module back into the LSTM. Figure 1
1445 illustrates the differences between these configurations; for convenience, we will use “post-
1446 gating” to refer to the latter mechanism and “pre-gating” to refer to the mechanism used
1447 in the simulations described in the main text. As noted in the *Discussion*, the primary con-
1448 sequence of having the gate on the output side is that the gate can be controlled based on
1449 information coming out of the hippocampus, in addition to all of the cortical regions that
1450 are used to control the gate in our pre-gating model. The post-gating mechanism has been
1451 more widely used in machine learning (Ritter et al., 2018; Ritter, 2019; Pritzel et al., 2017)
1452 because it is more powerful – since the gating function has access to activated episodic
1453 memories in the LCA, the model can close/open the gate depending on the properties of
1454 these activated memories.



1455 **Appendix 4 Figure 1.** Unrolled network diagrams for the pre-gating (A) versus the post-gating (B)
1456 models. The EM gate in the pre-gating model controls the degree to which stored memories are
1457 activated within the LCA module, but does not control the degree to which the activated memories
1458 are transmitted to the cortex. By contrast, the EM gate in the post-gating model controls the degree
1459 to which activated memories in the LCA module are transmitted to the cortex, but it does not control
1460 how these memory activations are computed in the first place.
1462

Since it is still unclear what kinds of episodic memory gating are implemented in the brain (see below for further discussion), we experimented with both mechanisms. We focused on the pre-gating model in the main text since it involves fewer assumptions – critically, it does not assume that the gating mechanism has access to the content of memories that are activated within the hippocampus. That said, the key results for the pre-gating model, reported in the main text, qualitatively hold for the post-gating model (Figure 2). In particular, the post-gating model also 1) retrieves much more from episodic memory in the DM condition, compared to the other two conditions (Figure 2A, B, C); 2) retrieves more when it is uncertain about the upcoming state (Figure 2D); 3) delays its recall time when the penalty is higher (Figure 2E); 4) adjusts its EM gate value as a function of the schema strength in a way that is similar to the pre-gating model (Figure 2F); and 5) shows the effect that midway-encoded memories hurt next-state prediction performance (Figure 2G, H – note that this also holds true when midway-encoded memories are present during meta-training). Importantly, while the aforementioned patterns replicate across the models, the results are not exactly the same – the retrieval policy for the post-gating model is often more flexible (i.e., it can adapt better to current conditions), since its EM gate can be controlled

1476

1477

1478

1479

1480

1481

1482

1483

1484

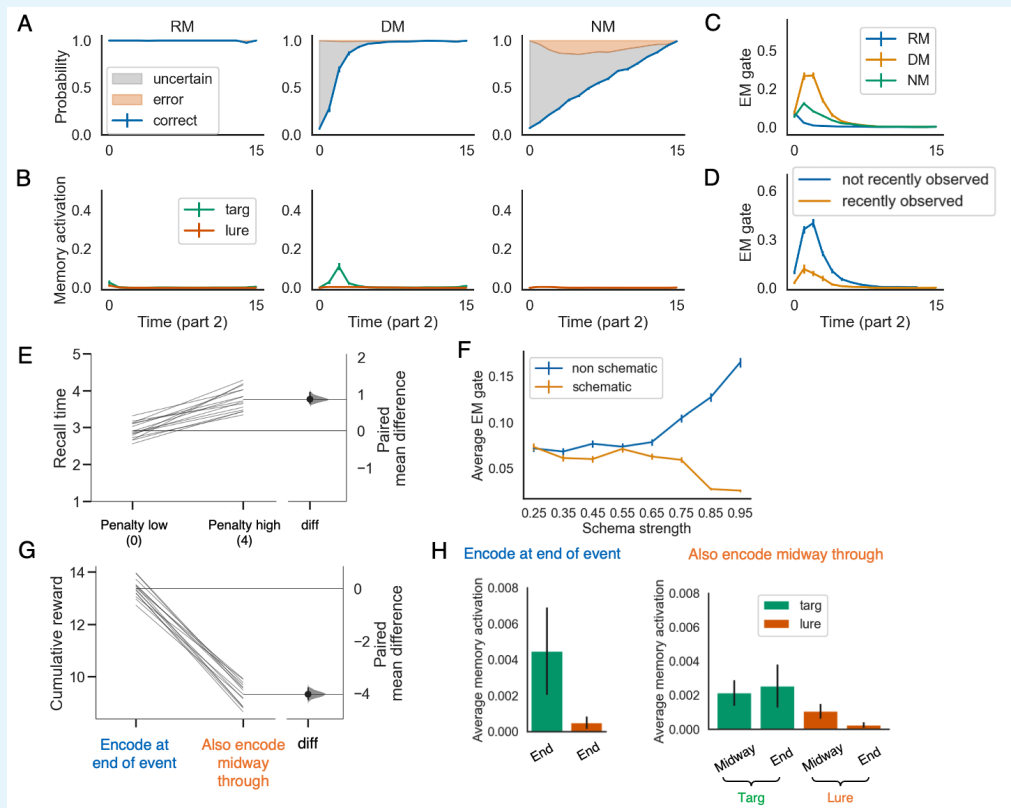
1485

1486

1487

1488

by the output of the EM module (in addition to the output of other cortical regions). For example, in the post-gating model, the EM gate layer of the cortical network is able to detect that relevant memories are not present in the NM condition, and it adapts to this by setting the EM gate to a lower value in the NM condition than the DM condition (Figure 2C) – that is, it learns to suppress retrieval when no memories are coming to mind. By contrast, the pre-gating model actually shows the opposite pattern – here, the EM gate layer can not detect the absence of relevant memories in the NM condition, but it *can* detect higher overall levels of uncertainty in the NM condition than the DM condition, which leads it to set the EM gate to a slightly higher value in the NM condition than the DM condition (see Figure 3C in the main text).



1489

1490

1491

1493

Appendix 4 Figure 2. The post-gating model qualitatively replicates key results obtained from the pre-gating model (compare to Figure 3, 4 in the main text). See text in this appendix for discussion. The errorbars indicate 1SE across 15 models.

One exciting future direction is to experimentally investigate how episodic memory gating works in the brain. The pre-gating and post-gating models make different predictions about the hippocampal activity: The post-gating model predicts that candidate episodic memory traces should be activated in the hippocampus at each time point; sometimes these activated traces are blocked (by the gate) from being transmitted to cortex, and sometimes they are allowed through. The pre-gating model predicts that activation of episodic memory traces in the hippocampus will be distributed more sparsely in time; on time points when the gate is closed, no activation should be transmitted from cortex to hippocampus, resulting in reduced activation of hippocampal memory traces (although there might be activation of these traces via recurrence within the hippocampus). Putting these points together, the pre-gating model appears to predict a large difference in hippocampal activation patterns as a function of whether the gate is closed or open; by contrast, the post-gating model appears to predict a smaller difference in hippocampal activation patterns as a func-

1502
1503
1504
1505
1506
1507
1508
1509
1510
1511
1512
1513
1514
1515
1516
1517
1518
1519
1520

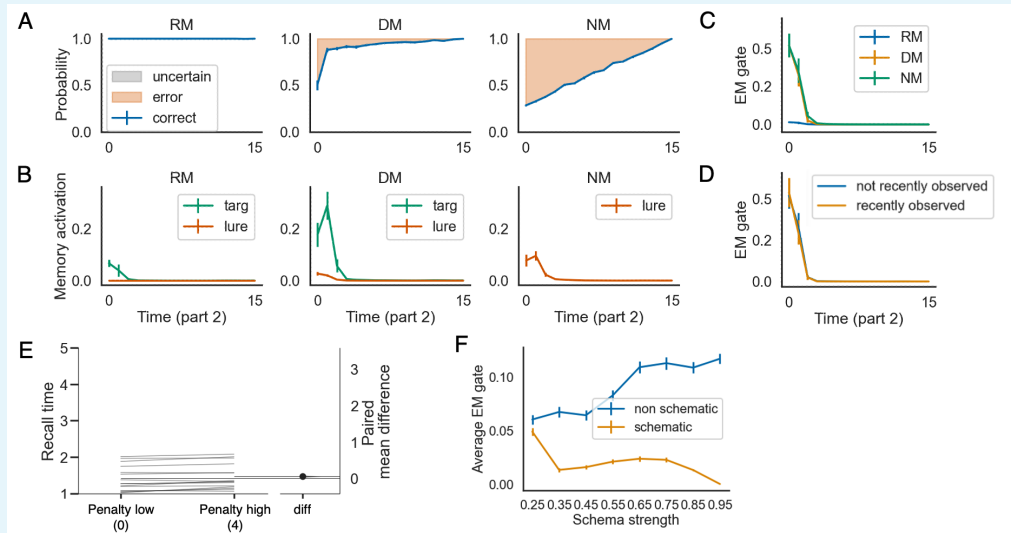
tion of whether the gate is closed or open.

However, this logic is complicated by the fact that the hippocampus is connected in a recurrent “big loop” with cortex (*Schapiro et al., 2017; Kumaran and Maguire, 2007; van Strien et al., 2009; Koster et al., 2018*) – in the post-gating model, even if the inputs to the hippocampus are the same when the gate is open vs. closed, the outputs to cortex will be different, which in turn will affect the inputs (from cortex) that hippocampus receives on the next time point. Thus, we would eventually expect differences in hippocampal activation in these conditions, even in the post-gate model. This suggests that, while it may be challenging to empirically tease apart the pre-gating and post-gating models, time-resolved methods like ECoG that can (in principle) distinguish between the “initial wave” of activity hitting the hippocampus after a stimulus and subsequent (recurrent) waves of activity would be most useful for this purpose. We should also note that the pre-gating and post-gating mechanisms are not mutually exclusive and it is possible that the brain deploys both of them.

1521 **Appendix 5**

1522

Training the model without reinforcement learning



1523

1524

1525

1526

1528

Appendix 5 Figure 1. Results from a “no-RL” model that was trained in an entirely supervised fashion, without reinforcement learning and without the option of giving a “don’t know” response – compare to Figure 3 in the main text; see text in this appendix for discussion. The errorbars indicate 1SE across 15 models.

1529

1530

1531

1532

1533

1534

1535

1536

1537

1538

1539

1540

1541

1542

1543

1544

1545

In the simulations shown in the main text, we trained the model using reinforcement learning (after supervised pre-training) and gave the model the option of responding “don’t know”, in which case it received no penalty or reward (see *Model training and testing* section above for details). Here, in Figure 1, we report the results from a model variant in which the model was trained in an entirely supervised fashion, without the option of responding “don’t know” – on each time point, the model was forced to predict the next state, and weights were adjusted based on the discrepancy between the predicted and actual states.

There are two important observations to make based on the results in Figure 1. The first observation is that the model is much less patient (i.e., it retrieves much earlier in part 2) when we take away the option of giving a “don’t know” response. This impatience can be seen by comparing the early time points of Figure 1C to the early time points of Figure 3C in the main text – EM gate values are much higher at early time points in the no-RL model. It can also be seen by comparing Figure 1E to Figure 3E in the main text – the average time-to-recall is much lower in the no-RL model. These findings confirm our claim (made in the main text) that the “don’t know” response makes the strategy of waiting to retrieve more viable, by allowing the model to escape being penalized on trials when it is waiting to retrieve from episodic memory.

The second observation is that, even without the option of responding “don’t know”, the learned retrieval policy of the no-RL model is still sensitive to certainty. This is shown in Figure 1B and C: Just like the model in the main text, the no-RL model recalls less information in the RM condition (when it is more certain about what will happen next) vs. the DM condition. The lack of a difference in EM gate value between “recently observed” and “not recently observed” features in Figure 1E suggests that the no-RL model might *not* be sensitive to certainty, but this is an artifact of the no-RL model’s impatience – the EM gate value is very high for early time points in both conditions, making it harder to observe a difference between conditions; in other simulations (not shown here) where we used a stronger

1551

1552

1553

1554

1555

1556

1557

1558

1559

1560

1561

1562

1563

1564

penalty manipulation to disincentivize early retrieval, the difference in recall levels between “recently observed” and “not recently observed” features was clearly visible in the no-RL model, reaffirming its sensitivity to certainty.

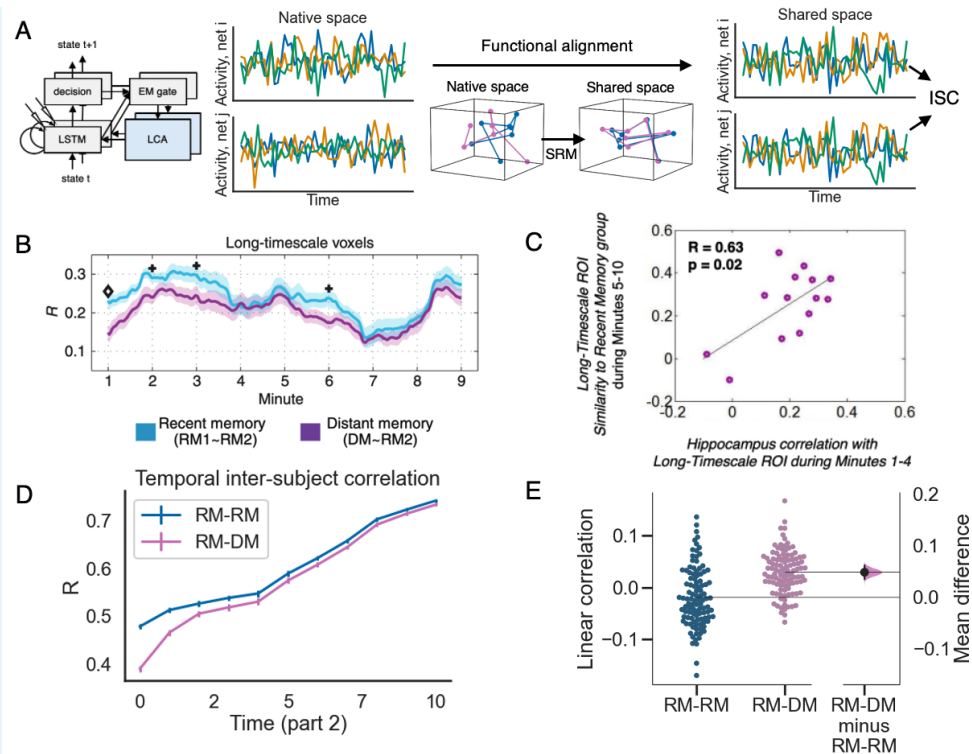
Taken together, the results from the no-RL model are very useful in clarifying what, exactly, is gained from the use of RL training with a “don’t know” option. In particular: having a “don’t know” response does not *cause* the model to have qualitatively distinct neural states as a function of certainty – these differences (described in Appendix 1 above) exist regardless of “don’t know” training, and can be used by the no-RL model to modulate its retrieval policy. Rather, the effect of RL training with the “don’t know” response is to make the model more patient, by giving it the option of waiting without penalty when it is uncertain.

1565 Appendix 6

1566 **Simulating inter-subject correlation results from Chen et al. (2016)**

1567 As discussed in the main text, *Chen et al. (2016)* found strong hippocampal-cortical activity
1568 coupling measured using inter-subject functional connectivity (ISFC; *Simony et al. 2016*) for
1569 DM participants, while the level of coupling was much weaker for participants in the RM
1570 and NM conditions (*Chen et al., 2016*). Here, we address some additional findings from this
1571 study that used temporal inter-subject correlation (ISC) as a dependent measure; temporal
1572 ISC tracks the degree to which the fMRI time series in a particular brain region is correlated
1573 across participants (*Hasson et al. 2004; Chen et al. 2016; Nastase et al. 2019*). Specifically,
1574 *Chen et al. (2016)* found that – at the start of part 2 – temporal ISC in DMN regions was lower
1575 between participants in the DM and RM conditions than between RM participants, suggest-
1576 ing differences in how DM and RM participants were interpreting the story; however, this
1577 gap in ISC decreased over the course of part 2, suggesting that these differences in inter-
1578 pretation between DM and RM participants decrease over time (Figure 1B). Furthermore,
1579 across participants, the degree to which the gap in ISC narrowed during the second half of
1580 part 2 was correlated with the amount of hippocampal-cortical activity coupling at the start
1581 of part 2 (Figure 1C; *Chen et al. 2016*). Taken together, these findings can be interpreted
1582 as showing that hippocampus is consulted more (as evidenced by increased hippocampal-
1583 cortical coupling) in the DM condition (where there are gaps in the situation model at the
1584 start of part 2) than the RM condition (where the situation model is more complete); the
1585 effect of this increased consultation of the hippocampus is to “fill in the gaps” and align the
1586 interpretations of the DM and RM participants (as evidenced by DM-RM ISC rising to the
1587 level of RM-RM ISC).

1588 To simulate these results, we trained 30 neural networks, then we assigned half of them
1589 to the RM condition and half to the DM condition. Next, we performed the temporal ISC anal-
1590 ysis used in *Chen et al. (2016)* by treating hidden-unit activity patterns as multi-voxel brain
1591 patterns. An important technical note is that running ISC across networks requires some
1592 form of alignment (i.e., so the time series for corresponding parts of the networks can be cor-
1593 related). Human fMRI data are approximately aligned across subjects, since brain anatomy
1594 is highly similar across people. However, when many instances of the same neural network
1595 architecture are trained on the same data, they tend to acquire different neural represen-
1596 tations, even though they represent highly similar mathematical functions (*Li et al., 2015;*
1597 *Dauphin et al., 2014; Meng et al., 2018*). That is, the same input can evoke uncorrelated neu-
1598 ral responses across different networks, although they produce similar outputs. For our
1599 purpose, this means that directly correlating hidden-layer activity patterns across neural
1600 networks will underestimate the similarity of representations across networks. Therefore,
1601 to simulate effects involving (human) inter-subject analyses, we need a way to align neural
1602 networks.



1603
1604
1605
1606
1607
1608
1609
1610
1611
1612
1613
1614
1615
1616
1617
1618
1619
1620
1622

Appendix 6 Figure 1. A) Illustration of how we computed inter-subject correlation (ISC) in the model (see text for details). B and C show the empirical results from [Chen et al. \(2016\)](#) (reprinted with permission) and D and E show model results. B) The sliding-window temporal inter-subject correlation (ISC) over time, during part 2 of the movie. The recent memory ISC, or RM-RM ISC, was computed as the average ISC value between two non-overlapping subgroups of the RM participants. The distant memory ISC, or RM-DM ISC, was computed as the average ISC between one sub-group of RM participants and the DM participants. Initially, the RM-DM ISC was lower than RM-RM ISC, but as the movie unfolded, RM-DM ISC rose to the level of RM-RM ISC. C) For the DM participants, the level of hippocampal-cortical inter-subject functional connectivity at the beginning of part 2 of the movie (minutes 1-4) was correlated with the level of RM-DM ISC later on (minutes 5-10). D) Sliding window temporal ISC in part 2 between the RM models (RM-RM) compared to ISC between the RM and DM models (RM-DM). The convergence between RM-DM ISC and RM-RM ISC shows that activity dynamics in the DM and the RM models become more similar over time (compare to part B of this figure). The errorbars indicate 1SE across 15 models. E) The correlation in the model between memory activation at time t and the change in ISC from time t to time $t + 1$, for the first 10 time points in part 2. Each point is a subject-subject pair across the two conditions. The 95% bootstrap distribution on the side shows that the correlation between memory activation and the change in RM-DM ISC is significantly larger than the correlation between memory activation and the change in RM-RM ISC (see text for details).

1623
1624
1625
1626
1627
1628
1629
1630

To accomplish this goal, we used the shared response model (SRM) ([Lu et al., 2018](#)) – a functional alignment procedure commonly used for multi-subject neuroimaging data ([Chen et al., 2015b](#); [Haxby et al., 2011, 2020](#)). Intuitively, this method applies rigid body transformation to align different network activities into a common space. We have previously shown that neural networks with highly overlapping training experience can be aligned well with SRM ([Lu et al., 2018](#)). Here, we used the Brain Imaging Analysis Kit (BrainIAK) implementation of SRM ([Kumar et al., 2020a,b](#)) to align our trained networks before computing ISC (Figure 1A).

Our simulation results qualitatively capture the findings from [Chen et al. \(2016\)](#). During part 2, DM-RM ISC starts lower than RM-RM ISC, but as the event unfolds, they gradually converge (Figure 1D). Moreover, in the DM condition, the level of memory activation at time t is correlated with the increment in DM-RM ISC from time t to time $t + 1$ (Figure 1E). As a

1631

1632

1633

1634

1635

1636

1637

1638

1639

1640

1641

1642

1643

comparison point, in the RM condition (where the model is not relying on episodic retrieval to fill in gaps in the situation model), memory activation does not correlate with the change in (RM-RM) ISC. Collectively, these results establish that episodic retrieval accelerates the convergence between model activations in the DM and RM conditions.

More generally, this result shows that one can capture inter-subject results with computational models. Experiments using inter-subject analyses and natural stimuli are becoming increasingly popular (*Nastase et al., 2019; Sonkusare et al., 2019; Hamilton and Huth, 2018; Nastase et al., 2020*); our simulation results provide a proof-of-concept demonstration of how computational models of memory can engage with this literature.

1644 Appendix 7

1645 **Model parameters**

1646 We implemented the model in PyTorch (*Paszke et al., 2017, 2019*). The numbers of hidden
1647 units for the LSTM layer and the decision layer were 194 and 128, respectively. The level
1648 of competition in the LCA module was 0.8. The initial cell state of the LSTM was a random
1649 vector \sim isotropic Gaussian(0, .1).

1650 During the meta-training phase, we used the Adam optimizer (*Kingma and Ba, 2014*).
1651 The initial learning rate was $7e-4$. The learning rate decayed by 1/2 if the average prediction
1652 accuracy minus mistakes stayed within 0.1% from the previous best loss for 30 consecutive
1653 epochs. The minimal learning rate was $1e-8$. We used orthogonal weight initialization with
1654 gain of 1 (*Saxe et al., 2014*), and we used supervised initialization for 600 epochs to help the
1655 model develop useful representations (*Misra et al., 2017; Nagabandi et al., 2017*). During
1656 the supervised initialization phase, the model was trained to predict the upcoming state;
1657 episodic memory and the “don’t know” unit were turned-off during this phase. After the su-
1658 pervised initialization phase, the model was trained with A2C (*Mnih et al., 2016*) for another
1659 400 epochs. We used entropy regularization with weight of 0.1 to encourage exploration.
1660 For every epoch, the model was trained on 256 events.