

Research Articles: Behavioral/Cognitive

Memory reactivation during learning simultaneously promotes dentate gyrus/CA_{2,3} pattern differentiation and CA₁ memory integration

<https://doi.org/10.1523/JNEUROSCI.0394-20.2020>

Cite as: J. Neurosci 2020; 10.1523/JNEUROSCI.0394-20.2020

Received: 17 February 2020

Revised: 11 November 2020

Accepted: 17 November 2020

This Early Release article has been peer-reviewed and accepted, but has not been through the composition and copyediting processes. The final version may differ slightly in style or formatting and will contain links to any extended data.

Alerts: Sign up at www.jneurosci.org/alerts to receive customized email alerts when the fully formatted version of this article is published.

1

2

3

Memory reactivation during learning simultaneously promotes dentate

4

gyrus/CA_{2,3} pattern differentiation and CA₁ memory integration

5

6

Robert J. Molitor^{1,2}, Katherine R. Sherrill², Neal W Morton², Alexandra A. Miller³, &

7

Alison R. Preston^{1,2,3*}

8

¹Department of Psychology, The University of Texas at Austin, Austin, TX, 78712, USA

9

²Center for Learning and Memory, The University of Texas at Austin, Austin, TX, 78712,

10

USA

11

³Department of Neuroscience, The University of Texas at Austin, Austin, TX, 78712,

12

USA

13

14

*Correspondence: apreston@utexas.edu

15

Abbreviated title: Reactivation promotes differentiation, integration (50/50 characters)

Number of figures: 4

Abstract word count: 250/250

Introduction word count: 645/650

Discussion word count: 1498/1500

Acknowledgements: This work was supported by NIH grants R01 MH100121, T32 MH106454, F31 NS103458, F32 MH114869, and The University of Texas at Austin Biomedical Imaging Center pilot grant 11042016a. The authors thank Christine Coughlin, Susannah Cox, Holly Hodge, Rosa Muñoz, Sharon Noh, Athula Pudhiyidath, Shilpa Rajagopal, and Hannah Roome for their assistance with data collection.

Conflict of interest: The authors declare no competing financial interest.

16

17 **Abstract**

18 Events that overlap with previous experience may trigger reactivation of existing
19 memories. However, such reactivation may have different representational
20 consequences within the hippocampal circuit. Computational theories of hippocampal
21 function suggest that dentate gyrus and CA_{2,3} (DG/CA_{2,3}) are biased to differentiate
22 highly similar memories, whereas CA₁ may integrate related events by representing
23 them with overlapping neural codes. Here, we tested whether the formation of
24 differentiated or integrated representations in hippocampal subfields depends on the
25 strength of memory reactivation during learning. Human participants of both sexes
26 learned associations (AB pairs, either face-shape or scene-shape), and then underwent
27 fMRI scanning while they encoded overlapping associations (BC shape-object pairs).
28 Both before and after learning, participants were also scanned while viewing indirectly
29 related elements of the overlapping memories (A and C images) in isolation. We used
30 multivariate pattern analyses to measure reactivation of initial pair memories (A items)
31 during overlapping pair (BC) learning, as well as learning-related representational
32 change for indirectly related memory elements in hippocampal subfields. When prior
33 memories were strongly reactivated during overlapping pair encoding, DG/CA_{2,3} and
34 subiculum representations for indirectly related images (A and C) became less similar,
35 consistent with pattern differentiation. Simultaneously, memory reactivation during new
36 learning promoted integration in CA₁, where representations for indirectly related
37 memory elements became more similar after learning. Furthermore, memory
38 reactivation and subiculum representation predicted faster and more accurate inference
39 (AC) decisions. These data show that reactivation of related memories during new

40 learning leads to dissociable coding strategies in hippocampal subfields, in line with
41 computational theories.

42

43

44

45

46 **Significance Statement**

47 The flexibility of episodic memory allows us to remember both the details that
48 differentiate similar events and the commonalities among them. Here, we tested how
49 reactivation of past experience during new learning promotes formation of neural
50 representations that might serve these two memory functions. We found that memory
51 reactivation during learning promoted formation of differentiated representations for
52 overlapping memories in the dentate gyrus/CA_{2,3} and subiculum subfields of the
53 hippocampus, while simultaneously leading to the formation of integrated
54 representations of related events in subfield CA₁. Furthermore, memory reactivation and
55 subiculum representation predicted success when inferring indirect relationships among
56 events. These findings indicate that memory reactivation is an important learning signal
57 that influences how overlapping events are represented within the hippocampal circuit.

58

59

60 **Introduction**

61 The hippocampus is composed of multiple subfields that contribute to memory
62 processing and representation. Computational models propose that the anatomical
63 properties of dentate gyrus and CA_{2,3} (DG/CA_{2,3}) make these subfields ideal for pattern
64 separation, or the automatic orthogonalization of highly similar cortical inputs through
65 sparse firing (Marr, 1971; Schapiro et al., 2017). In contrast, the characteristics of CA₁
66 have been proposed to mediate memory integration, or the formation of overlapping
67 representations that code the common features across related episodes (Eichenbaum
68 et al., 1999; Schlichting and Preston, 2015; Schapiro et al., 2017). Electrophysiological
69 research evinces such representational dissociations among subfields: DG/CA_{2,3}
70 ensembles elicit distinct firing patterns with only small changes in the perceptual
71 features of an environment, whereas CA₁ activity patterns change gradually as
72 environments become perceptually distinct (Leutgeb et al., 2004, 2007). Parallel work in
73 humans has shown that changes in DG/CA_{2,3} activation distinguish between highly
74 similar object images or objects that share a similar context, whereas CA₁ responses do
75 not (Bakker et al., 2008; Lacy et al., 2011; Dimsdale-Zucker et al., 2018). Subiculum,
76 the output structure of the hippocampal circuit (O'Mara et al., 2001), may contribute to
77 both pattern differentiation (Potvin et al., 2009) and integration (Schapiro et al., 2012).

78 However, such prior work has not considered how memory reactivation drives
79 dissociable representational strategies within hippocampus, allowing representation
80 learning to go beyond a simple transformation between external sensory input and
81 memory output. Classic computational learning models propose that memory
82 representations should adjust to predict likely outcomes in response to environmental

83 cues, with integration occurring when stimuli predict the same outcome and
84 differentiation when stimuli predict distinct outcomes (Rumelhart et al., 1986). However,
85 recent fMRI findings indicate that differentiation can also occur when stimuli share a
86 common association or outcome (Schlichting et al., 2015; Favila et al., 2016;
87 Zeithamova et al., 2018). In those studies, hippocampal representations were more
88 distinct for stimuli that shared a common outcome than stimuli with different outcomes.
89 Such differentiation cannot be explained in terms of automatic separation of external
90 input through sparse coding in DG/CA_{2,3}; rather, a recent theoretical perspective
91 proposes that memory reactivation may account for how hippocampal representations
92 change in the face of event overlap (Ritvo et al., 2019).

93 According to this theory, optimal learning reduces competition among memories
94 through either differentiation or integration (Ritvo et al., 2019). Although sensory overlap
95 in the environment is certainly one factor that might drive formation of optimal
96 representations that reduce ambiguity (Leutgeb et al., 2004, 2007; Lacy et al., 2011;
97 Yassa and Stark, 2011), what may be more essential is how overlapping sensory input
98 drives reactivation of competing memories. Reactivated memories may be the “target”
99 of learning more so than the sensory features that elicited reactivation. Thus, in the
100 present study, we went beyond considering perceptual similarity as the sole driver of
101 hippocampal representations and tested whether the reactivation of related memories in
102 cortex during learning results in dissociable subfield coding. We hypothesized that
103 memory reactivation would be modulated by event similarity across learning (Vieweg et
104 al., 2015) and may thus be the key factor mediating the degree of representational
105 overlap observed for similar events in hippocampal subfields (Ritvo et al., 2019). We

106 also hypothesized that integration and differentiation would not be mutually exclusive
107 outcomes in response to memory reactivation, but that reactivation would instead lead
108 to the simultaneous formation of complementary differentiated and integrated
109 representations in DG/CA_{2,3} and CA₁.

110 To test these predictions, we parametrically manipulated perceptual similarity
111 between overlapping events in an associative inference task (**Fig. 1**). Participants
112 studied initial pairs and were scanned using high-resolution fMRI while learning
113 overlapping pairs. We tested memory for the learned pairs and inferred knowledge of
114 the indirect relationships across pairs, with inference performance serving as a
115 behavioral index of integration (Shohamy and Wagner, 2008; Zeithamova et al., 2012).
116 Critically, we quantified how memory reactivation during overlapping event learning
117 impacted hippocampal subfield representation.

118

119 **Materials and Methods**

120 **Participants.** Thirty-two right-handed individuals (15 females, aged 18—31 years,
121 mean = 21.5 years) participated after giving informed consent in accordance with a
122 protocol approved by the Institutional Review Board at the University of Texas at Austin.
123 Participants received \$25/hour in compensation. Data from six participants were
124 excluded from the analyses: two participants due to excessive head motion, one
125 participant who withdrew from the experiment, two participants who had incomplete
126 scanning sessions (the post-exposure and/or localizer phases were not scanned), and
127 one participant for image artifacts in the functional scans that precluded analysis of the
128 pre-exposure and localizer phases. The remaining participants (n = 26, 14 females)

129 were included in the analyses. We determined our final sample size based on related
130 studies that used similar paradigms and analytical approaches (Zeithamova et al., 2012;
131 Schlichting et al., 2015; Dimsdale-Zucker et al., 2018). Furthermore, this sample size
132 gave us an estimated statistical power of over 0.99 to detect an effect of visual similarity
133 on across-episode inference accuracy based on pilot data from a separate group of
134 participants (n = 30, 22 females, aged 18—22 years, mean = 18.9 years; repeated
135 measures ANOVA resulting in partial eta squared (η^2) = 0.280).

136

137 **Stimuli.** Stimuli were 58 unfamiliar faces (half male, half female, all Caucasian), 58
138 unfamiliar scenes (half natural, half manmade), 671 black shapes generated in
139 MATLAB (see *Visual similarity manipulation during new encoding* for more
140 information), and 74 novel objects (Hsu et al., 2014; Schlichting et al., 2015). A subset
141 of the stimuli was organized into 32 triads consisting of three items (A, B, C) that were
142 used in the associative inference task (**Fig. 1A**). The A items consisted of faces (16)
143 evenly split by gender, and scenes (16) evenly split by natural and manmade; all B
144 items were shapes (56); all C items were novel objects (32). Another subset of stimuli
145 (42 faces, 42 scenes, 42 objects, and 42 shapes) were used in the localizer task and
146 were not seen during the associative inference task. Assignment of stimuli to the triads
147 and localizer task was randomized across participants. Stimuli were presented using
148 Psychtoolbox in MATLAB (Brainard, 1997; Pelli, 1997; Kleiner et al., 2007).

149

150 **Task procedure.**

151 **Initial pair (AB) learning.** Participants learned the initial pairs (AB) across four study-
152 test blocks. During the study phase, each of the 32 initial pairs was presented for 3.5s
153 with a 0.5s inter-trial interval (ITI). The A item (face or scene) was always presented on
154 the left and the B item (shape) was always presented on the right. After studying all of
155 the pairs, participants were tested using a 3-alternative forced choice (3 AFC) test.
156 Participants were cued with the A item on the top of the screen and had to choose
157 between the appropriate B item and two foils. The foils were shapes from other triads,
158 such that participants could not base their decision on the familiarity of the shapes.
159 Participants had 10s to respond on each trial. After the participant's response,
160 corrective feedback was provided at the end of each trial for 1s. Test trials were
161 separated by 0.5s ITI. Anatomical images were collected during this phase.

162
163 **Visual similarity manipulation during new encoding.** To examine how the similarity
164 of event elements affects memory reactivation and behavior, the visual similarity of the
165 linking element (the shape, or B item) in the associative inference task was
166 parametrically manipulated (**Fig. 1B**). We manipulated visual similarity based on prior
167 work showing that hippocampal subfield responses are modulated by visual feature
168 overlap among events (Bakker et al., 2008; Leutgeb et al. 2004, 2007; Lacy et al.,
169 2011). There was a total of four conditions: exact match, high similarity, low similarity,
170 and new. In the exact match condition, participants saw the exact same linking B shape
171 when learning the initial pairs (AB) and overlapping pairs (BC). In the high and low
172 similarity conditions, each shape seen in the overlapping pairs was a parametric morph
173 of a shape from one of the initial pairs. "Parent" shapes were generated by taking 16

174 points distributed along the perimeter of a circle, randomly translating each point, and
175 then connecting adjacent points to create edges using spline interpolation. The shapes
176 in the high and low similarity conditions were generated by taking two parent shapes
177 and averaging the coordinates of corresponding vertices using different weights. The
178 high similarity shapes were weighted 80% to one parent shape and 20% to the other
179 parent, while the low similarity shapes were weighted 70% to one parent and 30% to the
180 other. In the new condition, participants saw a new shape paired with a novel object,
181 making these pairs non-overlapping with the initial pairs. The new pairs thus served as
182 a baseline for associative learning. Each participant studied eight triads per visual
183 similarity condition.

184 Differences in subjective similarity between the high and low similarity items were
185 confirmed in an independent sample of nine participants (8 females, aged 18–22
186 years, mean = 19.4 years). Participants in this sample rated visual similarity between
187 parent shapes and shape morphs presented side by side using a 5-point Likert scale (1
188 = not at all similar, 5 = very similar) across 180 trials. Exact matches were rated as
189 more similar than high similarity morphs [$t_{(8)} = 6.255$, $p < 0.001$, Cohen's $d = 2.085$],
190 high similarity morphs were rated as more similar than low similarity morphs [$t_{(8)} =$
191 9.312 , $p < 0.001$, $d = 3.104$], and low similarity morphs were rated as more similar than
192 new items [$t_{(8)} = 10.021$, $p < 0.001$, $d = 3.340$]. One caveat to quantifying subjective
193 similarity using this approach is that the comparison does not involve a memory
194 component. It is possible that if we inserted a delay between the presentation of two
195 shapes, the observed subjective similarity function (**Fig. 1C**) may have differed; for
196 instance, the subjective similarity differences between the high and low similarity

197 conditions might have been less pronounced. While this measurement caveat might
198 influence interpretation of the subjective similarity judgments themselves, it has less
199 impact on interpretation of our central behavioral and neural analyses. We observe
200 differences in memory performance and reactivation between the similarity conditions
201 (including the high and low conditions) that indicate the four similarity conditions
202 differentially impacted processing (see **Results**). Furthermore, our neural analyses
203 assessing learning-related representational change focus on the high similarity
204 condition only and do not rely on comparisons to the other similarity conditions (see
205 ***Exposure of individual items before and after learning***).

206

207 ***Overlapping pair (BC) learning.*** After participants learned the initial pairs, they were
208 scanned while learning the overlapping pairs. This phase again consisted of four study-
209 test blocks. During the study phase, the 32 pairs were presented using an event related
210 design, with pairs presented for 3.5s followed by 8.5s ITI of fixation. The C item (object)
211 was always presented on the left and the B item (shape) was always presented on the
212 right. After each study phase, participants were tested on the BC pairs using a 3 AFC
213 test which was not scanned. Participants were cued with the C item on the top of the
214 screen and had to choose between the appropriate B item and two foils. Feedback was
215 not given during this phase. Participants had 10s to respond on each test trial, and trials
216 were separated by 0.5s ITI.

217

218 ***Exposure of individual items before and after learning.*** Before learning the initial
219 pairs and after learning the overlapping pairs, participants were exposed to individually

220 presented A and C items (faces, scenes, and objects) from the high similarity condition.
221 These exposure phases were limited to a single visual similarity condition to maximize
222 the number of presentations for each stimulus and improve estimation of task-related
223 activation patterns (see *Estimation of individual stimulus patterns before and after*
224 *learning*). Using a single similarity condition also allowed us to control for the effects of
225 visual similarity when calculating representational change. The high visual similarity
226 condition was used because prior work in humans has shown that highly visually similar
227 stimuli elicit differential responses in DG/CA_{2,3} and CA₁ (Lacy et al., 2011).

228 In each exposure run, participants were scanned while items were presented for
229 1s with a 3s ITI. While each item was on the screen, participants completed a change-
230 detection task by indicating via button press whether a superimposed black cross
231 changed color to green or blue 100ms to 200ms after stimulus onset (Kriegeskorte et
232 al., 2008; Schlichting et al., 2015). There were four repetitions of each item per run, and
233 a total of four runs each in the pre-exposure and post-exposure phases. Trials were
234 pseudorandomized such that items within a triad were presented with at least two
235 interleaved items from other triads. Additionally, 20% of trials were null (i.e., there was
236 no object or change detection task) to improve item-level activation estimation in the
237 analysis; these null trials were placed randomly between item presentation trials. Trial
238 order and timing was identical in the pre- and post-exposure phases. Accuracy on the
239 change detection task was monitored to ensure that participants were paying attention
240 to the task but was not considered further.

241 There was also a non-scanned pre-exposure phase for items from the exact
242 match, low similarity, and new conditions that occurred before the first scanned pre-

243 exposure run. The purpose of this phase was to equate familiarity of the A and C items
244 in the exact match, low similarity, and new conditions to items in the high similarity
245 condition prior to pair learning. The non-scanned exposure was similar to the scanned
246 exposure phases, except the ITI was 0.5s and there were no null trials.

247

248 **Associative inference (AC) test.** Following the post-exposure phase, participants were
249 given a surprise test on the indirect relationship between the A and C items that shared
250 a common associate (B). The inference test was performed inside the scanner but was
251 not scanned. In this phase, participants were cued with the C item (object) and could
252 choose between A items of the same category (i.e., 3 faces or 3 scenes). On face trials,
253 participants were instructed to choose the person who would most likely own the cued
254 object. On scene trials, they were instructed to choose the location in which they would
255 most likely find the cued object. Critically, at no point were participants explicitly
256 instructed about the visual similarity manipulation or the overlap across learning.
257 Participants were given 10s to respond. No feedback was given.

258

259 **Localizer.** After the inference test, participants were scanned in a localizer task. In this
260 task, participants viewed a series of stimuli from the four stimulus categories used in the
261 experiment: faces, scenes, shapes, and objects. Stimuli were presented in a blocked
262 design, with each block consisting of eight images presented for 2.5s each with 0.5s ITI.
263 During each stimulus block, participants completed a one-back memory task in which
264 they had to detect a repeated stimulus. There was one repeated stimulus in each block.
265 Accuracy on the one-back task was monitored to ensure that participants were paying

266 attention to the task but was not considered further. Blocks were separated by 8s of
267 fixation. Participants completed three runs of the localizer task, with two blocks per
268 stimulus type per run.

269

270 **fMRI data collection and preprocessing.** Data were collected with a 3T Siemens
271 Skyra. There was a total of 15 functional scans (TR = 2000ms, TE = 30ms, flip angle =
272 73°, 1.7mm isotropic voxels, echoplanar imaging, multiband acceleration factor = 3)
273 across the pre-exposure, overlapping pair study, post-exposure, and localizer phases.
274 Three field maps (TR = 589ms, TE = 5ms/7.46ms, 1.5×1.5×2mm voxels, flip angle = 5°)
275 were collected to correct for distortions in the magnetic field: one immediately before the
276 pre-exposure phase to correct the pre-exposure scans, one before the overlapping pair
277 study phase to correct the study and post-exposure scans, and one before the localizer
278 phase to correct the localizer scans. A T1-weighted 3D MPRAGE volume was collected
279 (TR = 1900ms, TE = 2.43ms, flip angle = 9°, 1mm isotropic voxels) to facilitate
280 alignment and normalization of the functional data to an anatomical template. Two
281 coronal T2-weighted structural scans, aligned perpendicular to the hippocampal long-
282 axis, were collected (TR = 13150ms, TE = 82ms, 0.4mm×0.4mm in-plane, 1.5mm thru-
283 plane) and then averaged for subfield segmentation.

284 Functional and anatomical images were preprocessed using FMRIB Software
285 Library version 5.0.9 (FSL: <http://fsl.fmrib.ox.ac.uk/fsl/>) and Advanced Normalization
286 Tools (ANTS) 2.1 (Avants et al., 2011). Functional scans were motion corrected using
287 MCFLIRT in FSL and then registered to the final overlapping pair study run using affine
288 transformations in ANTS. Non-brain structures were removed from the functional scans

289 and MPAGE using BET in FSL. Additional data processing was carried out using
290 FEAT (FMRI Expert Analysis Tool) Version 6.00, part of FSL. The following pre-
291 statistics processing was applied to all functional images; co-registration with the
292 MPAGE and field map-based EPI unwarping using FUGUE (Jenkinson, 2003); grand-
293 mean intensity normalization of the entire 4D dataset by a single multiplicative factor;
294 highpass temporal filtering (Gaussian-weighted least-squares straight line fitting, with
295 sigma = 64s). Spatial smoothing using a Gaussian kernel of FWHM 4mm was applied to
296 the overlapping pair learning and localizer scans.

297

298 **Regions of interest.** Anatomical regions of interest included whole-brain gray matter
299 for the reactivation analysis and hippocampal subfields for the neural coding analysis. A
300 whole-brain gray matter mask was created for each participant in native space using
301 FAST (Zhang et al., 2001), part of FSL, with the MPAGE. Gray matter masks were
302 then moved into functional resolution using linear transformations in ANTS. Within
303 hippocampus, activation patterns in subfields CA₁, a combined DG/CA_{2,3} region, and
304 subiculum were analyzed. Hippocampal subfields were identified in the head and body
305 of the hippocampus in native space by reverse normalizing masks from an open source
306 template with segmented subfields (Schlichting et al., 2019) to the average T2 coronal
307 image of each participant using non-linear SyN transformations in ANTS. This
308 procedure has been shown to provide results comparable to manual tracing (Schlichting
309 et al., 2019). Masks were then inspected and edited manually for each participant to
310 remove voxels outside the hippocampus and ensure accurate segmentation based on
311 established protocols (West and Gundersen, 1990; Duvernoy, 1998; Mai et al., 2007).

312 Finally, the subfield masks were transformed to the space of the functional scans by first
313 registering the average coronal image to the MPRAGE using linear transformations and
314 then applying the previously calculated transform to functional space.

315

316 **Quantification and statistical analysis.**

317 ***Decoding memory reactivation during overlapping event learning.*** To measure
318 reactivation of encoding patterns related to the initial pairs during overlapping pair
319 learning, we used a pattern classification analysis in PyMVPA (Hanke et al., 2009). If
320 participants reactivated related information (i.e., A face and scene items from AB pairs)
321 when learning overlapping pairs (BC shape-object pairs), then a pattern classifier
322 trained on the localizer data should be sensitive to the category of information (either
323 face or scene) that is being reactivated (Polyn et al., 2005; Kuhl et al., 2011;
324 Zeithamova et al., 2012). Thus, we trained the pattern classifier with data from the
325 localizer phase and then applied the classifier to the overlapping pair learning phase.
326 We operationalized memory reactivation as classifier evidence for the category of the A
327 items (i.e., face or scene) from the initial AB pairs related to the overlapping BC pairs.

328 We measured memory reactivation using a multi-step procedure. First, we ran a
329 whole-brain searchlight (Kriegeskorte et al., 2006) to identify regions where information
330 about A items was reinstated during overlapping pair learning. In each searchlight
331 sphere (radius = 3 voxels, volume = 123 voxels), a linear support vector machine was
332 trained to differentiate neural patterns from the localizer phase associated with faces,
333 scenes, objects, and shapes. To account for hemodynamic lag, each functional image
334 was labeled by taking the trial labels and time-shifting them forward by 4s (two TRs).

335 The trained classifier was then applied to neural patterns from the overlapping pair
336 learning phase, which was also time-shifted by 4s. Trial-level reactivation estimates
337 were extracted by taking classifier evidence for the category associated with the A item
338 of each triad (e.g., classifier evidence for faces for face-shape-object triads) for the two
339 TRs corresponding to the presentation of each pair. Classifier evidence values were
340 sorted into two sets: a reactivation set and baseline set. The reactivation set contained
341 classifier evidence values from the exact match, high similarity, and low similarity trials.
342 The baseline set contained face and scene evidence values from trials in the new
343 condition. Because shape-object pairs in the new condition did not overlap with any of
344 the previously learned pairs, they should not elicit reactivation of face or scene
345 memories. The final reactivation index was calculated in each sphere by taking the
346 difference between the average evidence for the reactivation set and the average of the
347 baseline set.

348 To test the significance of this reactivation index, we compared the actual
349 reactivation index to a null distribution in each searchlight sphere. The null distribution
350 was created over 1,000 iterations by shuffling classifier evidence values across the
351 reactivation and baseline sets and then re-calculating the reactivation index every
352 iteration. The center voxel of each searchlight sphere reported the proportion of the null
353 distribution with reactivation indices greater than or equal to the observed reactivation
354 index (i.e., p -value). To identify reactivation regions across participants, individual
355 participant searchlight maps were normalized to a group template for significance
356 testing. The p -value images were converted to z-statistic images and then warped to the
357 MNI 152 anatomical template (resampled to the resolution of the functional scans,

358 1.7mm isotropic voxels) using non-linear SyN transformations in ANTS. Voxel-wise,
359 nonparametric permutation testing was done using Randomise in FSL over 5,000
360 iterations (Winkler et al., 2014). Significant clusters were identified by applying a voxel
361 threshold of $p < 0.01$ (uncorrected) and a cluster threshold of $p < 0.05$. Thresholds were
362 calculated using the AFNI (Cox, 1996) function 3dClustSim with smoothness estimates
363 derived from the study phase using 3dFWHMx based on the spatial AutoCorrelation
364 Function (ACF). Cluster extent was determined using two-sided thresholding with
365 second-nearest neighbor clustering.

366 To confirm that the reactivation measure was not driven by a single stimulus
367 category, we further interrogated searchlight clusters to test whether reactivation varied
368 with stimulus category (face or scene) of the A item in a *post hoc* analysis. The
369 significant clusters identified in the reactivation searchlight analysis were converted to
370 binary masks and reverse-normalized into native space using ANTS. Then, the
371 reactivation analysis was repeated in each functional region of interest for every
372 participant. We used repeated measures ANOVA with region and stimulus category as
373 factors to assess whether reactivation in each region differed as a function of stimulus
374 category.

375 While our initial searchlight analysis localized regions in which reactivation
376 occurred above baseline, we also ran an independent searchlight to identify regions
377 where reactivation strength varied with visual similarity. This searchlight used a similar
378 approach to the analysis measuring overall reactivation, but with an additional level that
379 compared classifier evidence for reactivation between the exact match condition and
380 the other similarity conditions (i.e., the high similarity condition and low similarity

381 condition combined). The effect of similarity was calculated in each sphere by taking the
382 difference between the average evidence for the exact match condition and the average
383 evidence for the high and low similarity conditions combined. This difference was then
384 compared to a null distribution in each searchlight sphere, which was created over
385 1,000 iterations by shuffling classifier evidence values across the exact match and
386 similarity morph conditions. Normalization to the group template, statistical testing, and
387 cluster correction were identical to the searchlight identifying reactivation above
388 baseline.

389

390 ***Estimation of individual stimulus patterns before and after learning.*** We derived
391 estimates of neural activation patterns elicited by each of the A (faces, scenes) and C
392 (novel 3D objects) stimuli from the pre-exposure and post-exposure phases using a
393 general linear model (GLM) with a least squares–separate (LS-S) approach (Mumford
394 et al., 2012) in the native space of each participant. Each of the 16 objects from the
395 scanned pre-exposure phase (i.e., the eight A items and eight C items from the high
396 similarity condition) were modeled iteratively in each run of the pre- and post-exposure
397 phases separately (Schlichting et al., 2015).

398 Object presentations were modeled as a 1s event, and the regressor for each
399 object included all four repetitions within a scanning run. Each of the 16 object
400 regressors was convolved with the canonical double gamma hemodynamic response
401 function. Temporal filtering was then applied. The GLMs included additional confound
402 regressors: motion parameters, their temporal derivatives, framewise displacement
403 (FD), and DVARS (Power et al., 2012; Schlichting and Preston, 2014; Schlichting et al.,

404 2015). Additional motion regressors were added for time points during which head
405 motion exceeded both 0.5mm for FD and 0.5% change in BOLD signal for DVARS
406 (Power et al., 2012). Beta images were generated for each A and C item for every pre-
407 and post-exposure run, totaling 128 statistics images per participant.

408

409 ***Quantifying learning-related changes in hippocampal subfield neural similarity.***

410 Pattern differentiation and memory integration in hippocampus were indexed using a
411 representational similarity analysis (Kriegeskorte et al., 2008) implemented in PyMVPA
412 (Hanke et al., 2009). Searchlights were run separately within anatomically defined
413 DG/CA_{2,3}, CA₁, and subiculum in the native space of each participant. Within each
414 searchlight sphere (radius = 2 voxels, volume = 33 voxels) (Schapiro et al., 2012),
415 similarity matrices were generated by calculating the pairwise Pearson's correlation
416 values for the 128 statistics images corresponding to the A and C items in the pre-
417 exposure and post-exposure runs, transformed to Fisher's z. Then, change in pattern
418 similarity due to learning was measured by subtracting the pre-exposure similarity
419 values from the post-exposure similarity values in corresponding cells.

420 After the change in pattern similarity (hereafter referred to as Δ) was calculated,
421 Δ values were sorted depending on whether the value was for a within-triad comparison
422 or an across-triad comparison. These two sets of values allowed us to determine how
423 representational change was influenced by event overlap due to learning (within-triad
424 comparison set) relative to a baseline that simply reflected repeated exposure without
425 event overlap (across-triad comparison set). Importantly, only Δ values that reflected

426 across-run correlations were used to reduce bias that could be introduced from
427 autocorrelation in the BOLD signal (Mumford et al., 2012).

428 To assess the effect of reactivation during learning on representational change,
429 the within-triad Δ values were further subdivided based on the strength of memory
430 reactivation during learning of the overlapping pairs. For each participant, reactivation
431 strength was calculated for every triad by taking the mean reactivation index across the
432 network of regions identified in the reactivation searchlight analysis (**Fig. 3A**), averaged
433 across study blocks. Triads were then divided into stronger reactivation triads and
434 weaker reactivation triads using a median split on the average reactivation values.
435 Thus, within-triad Δ comparisons were further sorted into a stronger reactivation within-
436 triad Δ set and a weaker reactivation within-triad Δ set in each searchlight sphere.
437 Finally, all Δ sets were averaged to create three summary values: average within-triad
438 similarity change for stronger reactivation triads ($\Delta_{\text{Within stronger}}$), average within-triad
439 similarity change for weaker reactivation triads ($\Delta_{\text{Within weaker}}$), and average across-triad
440 similarity change (Δ_{Across}). We compared these summary values to determine whether
441 neural coding varied as a function of reactivation strength.

442 Neural coding was assessed using four searchlight contrasts (Schlichting et al.,
443 2015) (**Fig. 4B**). Two analyses identified hippocampal voxels for which there was
444 memory integration or differentiation across all triads, regardless of reactivation
445 strength. Integration_{Overall} was calculated as $(\Delta_{\text{Within stronger}} - \Delta_{\text{Across}}) + (\Delta_{\text{Within weaker}} -$
446 $\Delta_{\text{Across}})$, reflecting greater within-triad than across-triad similarity after learning across all
447 degrees of reactivation. Differentiation_{Overall} was calculated as $(\Delta_{\text{Across}} - \Delta_{\text{Within stronger}}) +$
448 $(\Delta_{\text{Across}} - \Delta_{\text{Within weaker}})$, reflecting lesser within-triad than across-triad similarity across all

449 degrees of reactivation. Two additional analyses identified voxels for which neural
450 coding varied as a function of reactivation strength ($\text{Integration}_{\text{Reactivation}}$ and
451 $\text{Differentiation}_{\text{Reactivation}}$). The $\text{Integration}_{\text{Reactivation}}$ searchlight identified voxels for which
452 integration occurred to a greater extent for stronger reactivation triads.
453 $\text{Integration}_{\text{Reactivation}}$ was calculated as $(\Delta_{\text{Within stronger}} - \Delta_{\text{Within weaker}})$. In contrast, the
454 $\text{Differentiation}_{\text{Reactivation}}$ searchlight identified voxels for which differentiation occurred to a
455 greater extent for stronger reactivation triads. $\text{Differentiation}_{\text{Reactivation}}$ was calculated as
456 $(\Delta_{\text{Within weaker}} - \Delta_{\text{Within stronger}})$.

457 The significance of each of these calculations was determined by comparing the
458 computed similarity change values to a null distribution in each searchlight sphere. The
459 null distribution was created over 1,000 iterations by shuffling cells across the Δ_{Within}
460 stronger, $\Delta_{\text{Within weaker}}$, and Δ_{Across} sets and then re-calculating the statistic of interest each
461 iteration. The center voxel of each searchlight sphere reported the proportion of the null
462 distribution with values greater than or equal to the observed similarity change (i.e., p -
463 value). Significant clusters were identified using the same method as the reactivation
464 searchlights, except the z-statistic images were warped to a functional-resolution
465 hippocampal template rather than the re-sampled MNI template for the group-level
466 analysis. Normalized searchlight maps were then masked by each anatomical
467 hippocampal subfield template prior to cluster correction to ensure clusters were
468 exclusive to each hippocampal subfield.

469 *Post hoc* analyses further interrogated the direction of representational change
470 observed in each subfield identified from this searchlight analysis. An important caveat
471 to these *post hoc* analyses is that they are not completely unbiased because they

472 compare sets of voxels pre-selected to exhibit specific effects based on the searchlight
473 contrasts. Thus, our follow-up analyses did not directly compare the Δ_{Within} values for the
474 stronger and weaker reactivation items. Our *post hoc* analyses instead focused on the
475 magnitude of Δ_{Across} values to test whether there were global shifts in neural similarity
476 across the pre- and post-exposure phases, in addition to comparing Δ_{Within} values to
477 Δ_{Across} values to quantify the degree of learning-related integration and differentiation.

478 For these *post hoc* analyses, similarity change in DG/CA_{2,3}, CA₁, and subiculum
479 was calculated for each participant in native space. The searchlight clusters identified
480 by the group searchlight analyses were converted into masks and reverse-normalized
481 into each participant's native space using non-linear transformations in ANTS. For each
482 participant, the native space clusters were then dilated with FSL using a 3x3x3 mm box
483 as a kernel. To ensure that clusters were still restricted to their respective subfield when
484 converted to participant native space, each cluster was masked using anatomical
485 subfield masks defined for each individual participant. One participant had a CA₁ cluster
486 in native space without a sufficient number of voxels for representational similarity
487 analysis (< 10 voxels) and was excluded from subsequent analysis of this subfield. For
488 the remaining participants, we computed the average similarity change within each
489 cluster separately for triads associated with stronger reactivation during learning, those
490 associated with weaker reactivation during learning, and the across-triad baseline.

491

492 ***Quantifying the relationship between neural measures and behavior.*** The
493 relationship between behavior and our neural measures of reactivation and
494 representational change was assessed using a Linear Ballistic Accumulator (LBA)

495 model to fit performance on the inference test (Morton et al., In press). For each
496 participant and subfield (CA₁, DG/CA_{2,3}, and subiculum), we calculated the z-score of
497 similarity change between A and C items from pre- to post-learning (Δ) for each triad.
498 We also calculated the z-score of A item reactivation across triads for each participant.
499 We then used the LBA model to fit behavioral responses and response times during the
500 AC inference test, using similarity change and reactivation as predictors of variability
501 between triads. We used a multilevel Bayesian approach to estimate mean slope
502 parameters reflecting the relationship between the neural measures and AC inference
503 performance. Positive slopes for the Δ measures indicate larger similarity values
504 between A and C items after learning are associated with faster and more accurate
505 inference. Positive slopes for the reactivation measure indicate that greater reactivation
506 is associated with faster and more accurate inference.

507

508 **Model definition.** The LBA model (Brown and Heathcote, 2008) assumes that, on each
509 trial, the starting point k of each accumulator is drawn randomly from a uniform
510 distribution on the interval $[0, A]$. Each accumulator then follows a line with a slope of d
511 until it reaches the response threshold b . On each trial, the slope d of accumulator i is
512 drawn from a normal distribution with mean v_i and standard deviation s (here, fixed at
513 1). The time for an accumulator to reach the threshold is $(b - k)/d$. We modeled the
514 three-alternative forced-choice inference tests using three accumulators with mean drift
515 rates v_1 (for the correct response) and v_2 (for the other two responses).

516 As derived in the initial description of the LBA model (Brown and Heathcote,
517 2008), the probability density function (PDF) for accumulator i at time t is:

518

$$f_i(t) = \frac{1}{A} \left[-v_i \phi \left(\frac{b - A - tv_i}{ts} \right) + s \phi \left(\frac{b - A - tv_i}{ts} \right) + v_i \Phi \left(\frac{b - tv_i}{ts} \right) - s \Phi \left(\frac{b - tv_i}{ts} \right) \right]$$

519 Where ϕ and Φ are the probability density function and cumulative distribution

520 functions, respectively, of the standard normal distribution. The cumulative distribution

521 function (CDF) for accumulator i at time t is:

$$F_i(t) = 1 + \frac{b - A - tv_i}{A} \Phi \left(\frac{b - A - tv_i}{ts} \right) - \frac{b - tv_i}{A} \Phi \left(\frac{b - tv_i}{ts} \right) \\ + \frac{ts}{A} \phi \left(\frac{b - A - tv_i}{ts} \right) - \frac{ts}{A} \phi \left(\frac{b - tv_i}{ts} \right)$$

522 The PDF for accumulator i hitting the threshold first, at time t , is the probability of

523 accumulator i finishing at time t , conditional on the other accumulators not having

524 finished yet:

$$\text{PDF}_i(t) = f_i(t) \prod_{j \neq i} (1 - F_j(t))$$

525 Because drift rate d is drawn from a normal distribution, there is some probability of no

526 accumulators finishing. Following prior work (Brown and Heathcote, 2008; Annis et al.,

527 2017), we conditionalized on the probability of at least one accumulator having a

528 positive drift rate:

$$P(\text{resp}) = 1 - \prod_{i=1}^N \phi \left(-\frac{v_i}{s} \right)$$

529 Non-decision time (e.g., time to perceive the test stimuli) was modeled as a fixed time

530 interval τ . The probability of a correct response at time t was:

$$P(\text{correct}, t) = \frac{\text{PDF}_1(t - \tau)}{P(\text{resp})}$$

531 The probability of an incorrect response at time t was:

$$P(\text{incorrect}, t) = \frac{2\text{PDF}_2(t - \tau)}{P(\text{resp})}$$

532 The model was implemented in Python 3.7 using PsiReact 0.2 (Morton et al., In
 533 press). We used Bayesian sampling to estimate parameters, using the No U-Turn
 534 Sampler (NUTS) implemented in pyMC 3.9.2. We fixed $s = 1$ and $b = 8$ to improve
 535 stability of parameter estimates. An intercept drift rate parameter $\beta_{0,i}$ for correct
 536 responses was estimated for each participant i . We also estimated the drift rate of
 537 incorrect responses $v_{2,i}$ for each participant. We used the within-participant z-score of
 538 similarity change for each subfield (e.g., $z_{CA1,ij}$) and reactivation estimates ($z_{\text{React},ij}$) to
 539 predict the drift rate on each trial j . Trial-level variability in drift rate was modeled as a
 540 linear combination of the similarity change and reactivation z-scores. The correct item
 541 drift rate $v_{1,ij}$ for participant i , trial j was:

$$v_{1,ij} = \beta_{0,i} + \beta_{CA1,i}z_{CA1,ij} + \beta_{DG/CA2,3,i}z_{DG/CA2,3,ij} + \beta_{\text{Subiculum},i}z_{\text{Subiculum},ij} + \beta_{\text{React},i}z_{\text{React},ij}$$

542 The slope parameters (e.g., $\beta_{CA1,i}$) were estimated for each participant i . To improve
 543 robustness of estimates for the individual participant parameters, we defined them as
 544 being drawn from a group-level normal distribution. The prior distributions for
 545 parameters were:

$$\beta_{0,i} \sim \text{Normal}(0, 4)$$

$$\beta_{CA1,i} \sim \text{Normal}(\mu_{CA1}, \sigma_{CA1})$$

$$\beta_{DG/CA2,3,i} \sim \text{Normal}(\mu_{DG/CA2,3}, \sigma_{DG/CA2,3})$$

$$\beta_{\text{Subiculum},i} \sim \text{Normal}(\mu_{\text{Subiculum}}, \sigma_{\text{Subiculum}})$$

$$\beta_{\text{React},i} \sim \text{Normal}(\mu_{\text{React}}, \sigma_{\text{React}})$$

$$v_{2,i} \sim \text{Normal}(\mu_2, \sigma_2)$$

$$\tau \sim \text{Unif}(0, 2)$$

$$A \sim \text{Unif}(0, 8)$$

546 Prior distributions for group-level parameters were:

$$\mu_{CA1}, \mu_{DG/CA2,3}, \mu_{Subiculum}, \mu_{React}, \mu_2 \sim \text{Normal}(0, 4)$$

$$\sigma_{CA1}, \sigma_{DG/CA2,3}, \sigma_{Subiculum}, \sigma_{React}, \sigma_2 \sim \text{Gamma}(1.5, 0.75)$$

547 For each of 4 chains, there was a tuning phase of 1,000 iterations with a target
 548 acceptance rate of 0.99, followed by 5,000 samples. Convergence was assessed using
 549 bulk effective sample size and rank-normalized split potential scale reduction statistic \hat{R}
 550 (Vehtari et al., 2019). We assessed the fit of the model by calculating mean posterior
 551 parameters for each trial as well as simulating responses and response times. We
 552 simulated 50 replications of each trial to obtain a robust estimate of model performance.
 553 Finally, we calculated the 95% high-density interval for each of the group-level mean
 554 parameters (e.g., μ_{CA1} for CA_1) to determine whether they were different from zero,
 555 indicating a relationship between similarity change or reactivation and AC inference
 556 performance.

557

558 **Results**

559 **Behavioral performance.** By the end of the initial pair (AB) learning phase, participants
 560 had formed strong memories of the face-shape and scene-shape pairs. All participants
 561 were above chance on the final test (mean proportion correct = 0.91, standard deviation
 562 [SD] = 0.01) and were therefore included in subsequent analyses. Memory for the
 563 overlapping (BC) shape-object pairs was influenced by the visual similarity of the linking
 564 item across learning (**Fig. 2A, Fig. 2B**). A repeated measures ANOVA with the within-

565 subjects factors of overlapping pair block (1, 2, 3, 4) and visual similarity (exact match,
566 high similarity, low similarity, new) revealed that visual similarity modulated memory
567 accuracy [main effect of block, $F_{(3,75)} = 79.93$, $p < 0.001$, $\eta^2 = 0.762$; block \times visual
568 similarity interaction, $F_{(9,225)} = 2.88$, $p = 0.003$, $\eta^2 = 0.103$] and response time [main
569 effect of similarity on correct trials, $F_{(3,72)} = 5.14$, $p = 0.003$, $\eta^2 = 0.176$]. For the first
570 learning block of overlapping pairs, performance was superior (**Fig. 2A**) when the
571 linking item (B) was an exact match to the initially learned pairs (AB) relative to all other
572 conditions. There was an effect of visual similarity in the first test block [effect of visual
573 similarity in the first run, $F_{(3,75)} = 6.901$, $p < 0.001$, $\eta^2 = 0.216$] but not in subsequent
574 runs [F -values ≤ 0.479 , all $p \geq 0.698$, all $\eta^2 \leq 0.019$]. In the first run, *post hoc* paired
575 t -tests revealed that accuracy was highest for pairs with an exact match relative to all
576 other pairs [compared to high similarity: $t_{(25)} = 3.33$, $p = 0.003$, $d = 0.654$; low similarity:
577 $t_{(25)} = 4.52$, $p < 0.001$, $d = 0.894$; new: $t_{(25)} = 2.74$, $p = 0.011$, $d = 0.539$]. Performance
578 was greater for high similarity pairs than low similarity pairs [$t_{(25)} = 2.306$, $p = 0.03$, $d =$
579 0.459]. There was no difference in performance between the high similarity and new
580 pairs [$t_{(25)} = 0.87$, $p = 0.394$, $d = 0.172$] or the low similarity and new pairs [$t_{(25)} = 0.76$, p
581 $= 0.452$, $d = 0.151$]. When collapsed across block, pairs with exact matches had the
582 fastest response time (**Fig. 2B**) on correct trials [compared to all other conditions, t -
583 values ≥ 2.206 , all $p < 0.05$, all $d \geq 0.445$]. Response time did not differ between the
584 high similarity, low similarity, or new pairs [all t -values ≤ 1.748 , all $p > 0.05$, all $d \leq$
585 0.348].

586 Visual similarity of the linking item also influenced cross-episode inference
587 accuracy [$F_{(3,75)} = 26.61$, $p < 0.001$, $\eta^2 = 0.516$]. Participants were more likely to infer a

588 relationship among indirectly related memory elements (AC) when the linking item (B)
589 was an exact match or highly similar across overlapping pairs (**Fig. 2C**). Inference
590 performance did not differ between exact match and high similarity triads [$t_{(25)} = 1.20$, p
591 $= 0.24$, $d = 0.230$], but performance for exact match triads was superior to both low
592 similarity triads [$t_{(25)} = 6.82$, $p < 0.001$, $d = 1.327$] and new triads [$t_{(25)} = 6.61$, $p < 0.001$,
593 $d = 1.286$]. Likewise, performance for high similarity triads exceeded low similarity triads
594 [$t_{(25)} = 5.05$, $p < 0.001$, $d = 0.987$] and new triads [$t_{(25)} = 5.38$, $p < 0.001$, $d = 1.055$].
595 Inference did not differ between the low similarity and new triads [$t_{(25)} = 1.17$, $p = 0.254$,
596 $d = 0.224$]. However, performance for low similarity triads was reliably better than
597 chance [$t_{(25)} = 2.22$, $p = 0.04$, $d = 0.435$], whereas performance for new triads was not
598 [$t_{(25)} = 0.47$, $p = 0.64$, $d = 0.093$].

599 Inference decisions were also faster for the exact match and high similarity
600 conditions relative to the new (or non-overlapping) condition [$F_{(3,72)} = 11.79$, $p < 0.001$,
601 $\eta^2 = 0.329$], with inferences for the exact match condition being fastest overall (**Fig. 2D**).
602 Response time was faster for exact match triads relative to high similarity conditions
603 [$t_{(25)} = 3.41$, $p = 0.002$, $d = 0.669$] and new triads [$t_{(24)} = 5.00$, $p < 0.001$, $d = 0.999$], but
604 no different from low similarity triads [$t_{(25)} = 1.64$, $p = 0.114$, $d = 0.321$]. Response time
605 was faster for high similarity triads compared with new triads [$t_{(24)} = 2.93$, $p = 0.007$, $d =$
606 0.585], but did not differ from low similarity triads [$t_{(25)} = 1.11$, $p = 0.28$, $d = 0.217$]. Low
607 similarity triads were faster than new triads [$t_{(24)} = 3.86$, $p = 0.001$, $d = 0.773$]. Together,
608 these findings show that associative memory and cross-episode inference, two
609 processes that are thought to be supported by hippocampal subfields (Schapiro et al.,

610 2017), are influenced by the perceptual similarity of shared event elements, with
611 facilitated performance with higher levels of cross-episode similarity.

612

613 **Reactivation of overlapping memories during learning.** To test how cortical memory
614 reactivation during overlapping pair learning impacts hippocampal subfield
615 representations, we first used a searchlight analysis to identify where information about
616 the initial pairs was reactivated in cortex during learning. Within each searchlight
617 sphere, a pattern classifier was trained on data from a localizer phase and then applied
618 to the overlapping pair study phase (Zeithamova et al., 2012). The searchlight identified
619 regions in which classifier evidence for the target category of the related item (face or
620 scene A items from the initial pairs) exceeded a baseline index of classifier evidence for
621 the same category derived from the new (or non-overlapping) trials. We found evidence
622 that related memories were reactivated when learning the overlapping pairs in posterior
623 cingulate cortex, occipital cortex, and parietal cortex (**Fig. 3A**).

624 Importantly, there were no differences in reactivation strength as a function of A
625 item category (face, scene) across regions identified in the searchlight analysis (**Fig.**
626 **3B**). A repeated measures ANOVA with the within-subjects factors of region (left
627 parietal, right parietal, cingulate, superior occipital, inferior occipital) and stimulus
628 category (face, scene) demonstrated that reactivation varied across regions [main effect
629 of region, $F_{(4,100)} = 2.84$, $p = 0.028$, $\eta^2 = 0.102$] but did not differ by stimulus category
630 [main effect of category, $F_{(1,25)} = 0.002$, $p = 0.967$, $\eta^2 = 0$; category \times region interaction,
631 $F_{(4,100)} = 0.375$, $p = 0.826$, $\eta^2 = 0.015$]. Thus, our results were not driven by a single

632 stimulus category and reflect memory reactivation rather than the engagement of
633 category-specific processing regions.

634 We further tested whether visual similarity of the shared B item across learning
635 influenced the strength of memory reactivation for the A items. We predicted that
636 memory reactivation during learning would be stronger for pairs linked by a more
637 visually similar item. Using a similar approach to the previous analysis, a separate
638 searchlight analyses identified regions where classifier evidence for the related A item
639 was greater for the exact match condition than the high and low similarity conditions.
640 Consistent with our hypothesis, we found that the similarity of event components
641 modulated the strength of memory reactivation in left parietal cortex and occipital cortex
642 (**Fig. 3C**).

643

644 **Memory reactivation impacts neural coding in hippocampal subfields.** To test our
645 hypothesis that reactivation of related memories during new encoding would lead to
646 dissociable representation of overlapping memories in DG/CA_{2,3} and CA₁, we quantified
647 hippocampal subfield coding as a function of memory reactivation strength during
648 learning. Both before and after learning the pairs, participants were scanned while
649 viewing individual images of the A and C items from overlapping pairs in the high
650 similarity condition (**Fig. 1A**). We indexed differentiation and integration by measuring
651 learning-related changes in pattern similarity for indirectly related A and C items from
652 the same triad (Schlichting et al., 2015). Similarity changes within the same triad were
653 compared to a baseline of similarity changes between items in different triads. We
654 measured differentiation by testing for a decrease in pattern similarity between A and C

655 items after learning (**Fig. 4A**). In contrast, integration would be marked by increased
656 pattern similarity among indirectly related A and C items, reflecting formation of
657 overlapping codes for related memories (**Fig. 4A**).

658 To assess the impact of memory reactivation during learning on neural coding of
659 indirectly related memory elements, we calculated representational change for triads
660 based on the strength of reactivation across overlapping learning trials. For each
661 participant, we sorted overlapping pairs into those associated with stronger and weaker
662 reactivation of the corresponding initial pair, based on a median split of averaged
663 reactivation indices across all clusters identified in the reactivation searchlight (**Fig. 3A**).
664 We then compared neural coding between indirectly related A and C items associated
665 with different levels of reactivation. Critically, all analyses assessing representational
666 change in hippocampal subfields were based on data from high similarity triads only.
667 This approach holds the visual similarity of the linking item constant, providing a critical
668 test of whether memory reactivation mediates representational change above and
669 beyond alterations of the physical environment.

670 We ran four searchlight analyses within individual hippocampal subfields to test
671 for the effects of reactivation on learning-related representational change for indirectly
672 related memory elements (**Fig. 4B**). First, we used two searchlight analyses to identify
673 hippocampal regions that showed differentiation or integration of A and C items
674 regardless of the degree of memory reactivation during overlapping pair learning
675 ($\text{Differentiation}_{\text{Overall}}$ and $\text{Integration}_{\text{Overall}}$, respectively) and observed no significant
676 effects within hippocampus. Instead, we predicted that the representational similarity of
677 indirectly related items in hippocampal subfields would depend on the strength of

678 memory reactivation during learning of the overlapping pairs. To test this hypothesis, we
679 ran two additional searchlight analyses that looked for an interaction between learning-
680 related representational change and memory reactivation; these searchlights isolated
681 hippocampal regions showing either differentiation or integration on trials with stronger
682 reactivation during overlapping pair learning ($\text{Differentiation}_{\text{Reactivation}}$ and
683 $\text{Integration}_{\text{Reactivation}}$).

684 We found that stronger reactivation of initial pair memories during learning of the
685 overlapping pairs had different consequences on the direction of representational
686 change observed in hippocampal subfields. When initial (A) memories were strongly
687 reactivated during overlapping (BC) pair learning, DG/CA_{2,3} pattern similarity decreased
688 between A and C items from pre- to post-learning (**Fig. 4C, Fig. 4D;**
689 $\text{Differentiation}_{\text{Reactivation}}$). Subiculum exhibited the same pattern as DG/CA_{2,3}, with
690 stronger reactivation leading to decreased pattern similarity for A and C items. In
691 contrast, CA₁ showed an opposing pattern of representational change when memory
692 reactivation was stronger, with increased similarity among A and C items post-learning
693 (**Fig. 4C, Fig. 4D;** $\text{Integration}_{\text{Reactivation}}$). These findings suggest that representation of
694 overlapping memories within hippocampal subfields is contingent on memory
695 reactivation during learning, with the same conditions leading to dissociable
696 representational codes within DG/CA_{2,3}, CA₁, and subiculum.

697 Finally, we performed a series of *post hoc* analyses on each hippocampal
698 subfield identified in the searchlight analysis to further understand how reactivation
699 modulated coding in each region. We first quantified whether there were any global
700 shifts in neural similarity simply as a function of learning by calculating the across-triad

701 Δ for unrelated A and C items (i.e., the across-triad baseline). Across-triad Δ was not
702 significantly different from zero in CA₁ [$t_{(24)} = 0.383$, $p = 0.705$, $d = 0.077$] or subiculum
703 [$t_{(25)} = 1.233$, $p = 0.229$, $d = 0.242$], but was greater than zero for DG/CA_{2,3}, [$t_{(25)} =$
704 3.431 , $p = 0.002$, $d = 0.673$]. These results demonstrate the importance of comparing
705 similarity change for related events to a baseline, as even unrelated items may change
706 in similarity after learning.

707 Next, we compared the within-triad Δ for triads associated with strong
708 reactivation to the across-triad Δ baseline as a validation our searchlight results (**Fig.**
709 **4D**). As mentioned previously, a caveat to this analysis is that the results are potentially
710 biased by selecting voxels identified in the neural coding searchlight analysis.
711 Consistent with the predicted patterns of the searchlight contrasts (**Fig. 4B**), we found
712 evidence for differentiation, whereby neural similarity change for triads associated with
713 strong reactivation was less than the across-triad baseline in DG/CA_{2,3} [$t_{(25)} = 2.298$, $p =$
714 0.030 , $d = 0.451$] and subiculum [$t_{(25)} = 3.158$, $p = 0.004$, $d = 0.619$]. Within CA₁, we
715 showed a trend for integration with greater similarity within triads associated with
716 stronger reactivation post learning relative to the across-triad baseline [$t_{(24)} = 1.766$, $p =$
717 0.090 , $d = 0.353$]. Together, these *post hoc* analyses support the outcome of the
718 searchlight analysis and show that representation of overlapping events in subfields is
719 influenced by the reactivation of related memories during learning.

720 As an exploratory analysis, we also quantified within-triad Δ for triads associated
721 with weaker reactivation during learning. We found evidence for integration in DG/CA_{2,3}
722 [$t_{(25)} = 3.709$, $p = 0.001$, $d = 0.727$] and a trend in subiculum [$t_{(25)} = 1.849$, $p = 0.076$, $d =$
723 0.363], wherein Δ for triads associated with weaker reactivation was greater than that

724 observed for the across-triad baseline. This result suggests that representational shifts
725 in DG/CA_{2,3} may vary as a function of the level of competition, which may be different
726 when memories are strongly or weakly reactivated. No differences from baseline were
727 observed for triads associated with weaker reactivation in CA₁ [$t_{(24)} = 1.062$, $p = 0.299$, d
728 = 0.212].

729

730 **Memory integration supports inference decisions.** We used a Bayesian multilevel
731 model to examine the relationship between similarity change after learning (i.e.,
732 integration or differentiation) and performance on the AC inference test. We also
733 examined the relationship between reactivation of related memories during learning and
734 inference performance. One participant was excluded from this analysis due to an
735 insufficient number of voxels in CA₁ (< 10 voxels). We used an LBA model to
736 simultaneously model inference accuracy and response times. We used Bayesian
737 sampling with the model to estimate the slope of relationships between inference
738 performance and triad-level variability in similarity change and memory reactivation. We
739 first assessed whether the Bayesian sampling was converged. There were no
740 divergences during sampling; for each parameter in the model, \hat{R} was less than 1.00102
741 and the effective sample size was at least 5225. These results indicate that the
742 sampling successfully converged, and there were sufficient samples to estimate each
743 parameter.

744 We used mean posterior parameters to simulate model responses and found that
745 there was a good fit to the observed accuracy (**Fig. 5A**) and response times (**Fig. 5B**)
746 on the inference test, with the exception of a small number of trials with very long

747 response times. The mean slope parameters for learning-related change (**Fig. 5C**) were
748 positive for subiculum (95% high-density interval = [0.043, 0.477], $d = 1.37$) and
749 memory reactivation (HDI = [0.005, 0.437], $d = 1.51$). The slope parameters for CA₁
750 (HDI = [-0.189, 0.244], $d = 0.15$) and DG/CA_{2,3} (HDI = [-0.393, 0.102], $d = 0.50$) were
751 not different from zero. The 95% high-density intervals for the other model parameters
752 were: $A = [2.059, 5.601]$, $\tau = [0.00009, 0.515]$, $\mu_2 = [0.130, 0.812]$, $\sigma_2 = [0.191, 0.831]$,
753 $\sigma_{CA1} = [0.004, 0.458]$, $\sigma_{DG/CA2,3} = [0.010, 0.577]$, $\sigma_{Subiculum} = [0.002, 0.408]$, and
754 $\sigma_{React} = [0.0002, 0.327]$. These results indicate that greater memory reactivation during
755 learning and greater AC similarity after learning in subiculum predict faster and more
756 accurate inference at the level of individual trials.

757

758 **Discussion**

759 Our results indicate that reactivated memories guide how representations of related
760 events are organized within the hippocampal circuit. Reactivation of prior memories
761 during encoding of new, overlapping events predicted across-episode inference
762 performance and had different consequences for representation in hippocampal
763 subfields; strong reactivation led to differentiation of overlapping memories within
764 DG/CA_{2,3} and subiculum, while simultaneously promoting integration of those same
765 memories in CA₁. Prior work has focused on explaining hippocampal subfield coding in
766 terms of a transfer function through which changes in environmental cues lead to
767 differential neural output (Leutgeb et al., 2004, 2007; Lacy et al., 2011; Yassa and
768 Stark, 2011). Here, we show that changes in perceptual input are not the only factor
769 determining representation learning within hippocampal subfields. Rather, our data

770 indicate that hippocampal subfield coding is further driven by the degree to which a new
771 experience triggers reactivation of related episodes. Our results thus extend prior
772 findings to show—at a representational level—that cortical memory reactivation drives
773 dissociations in hippocampal subfield coding in the face of competition between highly
774 similar memories.

775 Prior work on hippocampal representation has primarily conceptualized subfield
776 coding as an automatic process in response to environmental changes, wherein
777 sensory inputs are assumed to be the main driver of hippocampal responses. For
778 instance, early electrophysiological studies in rodents measured how place field
779 responses in hippocampal subfields remapped as animals navigated environments with
780 gradually changing perceptual features (Guzowski et al., 2004; Lee et al., 2004;
781 Leutgeb et al., 2004, 2007; Vazdarjanova and Guzowski, 2004). Such work revealed
782 that small changes in environmental features led to dramatic changes in DG and CA₃
783 responses, reflecting orthogonalization of input patterns. In contrast, CA₁ responses
784 changed gradually, scaling linearly with the amount of perceptual change between
785 environments; for environments that were more perceptually similar, CA₁ activity showed
786 a greater overlap in responding. Prior work in humans took a similar approach,
787 presenting participants with pairs of highly similar visual images (e.g., pictures of two
788 different apples) and measuring the magnitude of hippocampal subfield responses to
789 both images (Bakker et al., 2008; Lacy et al., 2011). In those studies, DG/CA_{2,3} showed
790 a novelty response for both highly similar images, suggesting separate coding of the
791 two images. CA₁ and subiculum responses to the second, highly similar image from a

792 pair, however, were suppressed relative to the presentation of the first pair member,
793 suggesting similar representation of the paired images.

794 While past animal and human work has revealed important dissociations
795 between hippocampal subfield processing, our findings build upon that work to show
796 that hippocampal representation learning is not simply a passive process, but instead is
797 actively influenced by memory reactivation (Hulbert and Norman, 2015; Kim et al., 2017;
798 Ritvo et al., 2019). We show that hippocampal subfield dissociations are most apparent
799 when past memories are strongly reactivated, producing a competitive learning state
800 that promotes differentiation in DG/CA_{2,3} and subiculum, simultaneously with integration
801 in CA₁. Our data thus indicate the need to quantify both the perceptual similarity among
802 events and how overlapping perceptual features trigger memory reactivation to fully
803 account for how dissociable representations emerge within the hippocampal circuit. One
804 interesting aspect of the prior human work described above is that dissociations among
805 subfields depended on the nature of the task being performed (Kirwan and Stark, 2007;
806 Bakker et al., 2008; Lacy et al., 2011). When the critical experimental manipulation (i.e.,
807 the visual similarity among items) was incidental to the task participants performed,
808 dissociations between subfields were observed (Bakker et al., 2008; Lacy et al., 2011).
809 However, when the same stimuli and presentation procedures were combined with an
810 intentional task focus, dissociations were less apparent (Kirwan and Stark, 2007). The
811 mechanistic source of these divergent findings has yet to be revealed. By quantifying
812 memory reactivation during tasks with an intentional or incidental focus, further insights
813 might be gained about how task goals influence the dynamics of how memory
814 competition impacts neural representation (Richter et al., 2016).

815 Our findings may be conceptualized in terms of supervised and unsupervised
816 models of learning, which each focus on different learning targets. Whereas supervised
817 learning is directed by matching representations to sensory cues observed directly in
818 the environment, unsupervised learning adjusts representations to reduce competition
819 between a current experience and reactivated memory representations triggered by the
820 new event (Ritvo et al., 2019) through integration or differentiation. While learning likely
821 reflects a balance between supervised and unsupervised mechanisms, our findings
822 indicate that reactivated memories are an important facet of how dissociable coding
823 strategies emerge across hippocampal subfields.

824 To date, only one other study in humans has used multivariate representational
825 analyses to quantify a dissociation between hippocampal subfields, specifically when
826 individuals retrieved information about shared or distinct spatial contexts (Dimsdale-
827 Zucker et al., 2018). That study showed that items learned within the same spatial
828 context elicited overlapping activation patterns in CA₁ and differentiated patterns in
829 DG/CA_{2,3} during retrieval relative to items that did not share contextual information. The
830 present findings differ from that study in several key ways. First, the prior study
831 measured subfield codes during memory retrieval, while our work reveals the active
832 learning processes that drive formation of dissociable subfield representations.
833 Specifically, that prior study did not quantify how reactivation of similar memories, either
834 during learning or retrieval, related to hippocampal subfield coding. Here, we show a
835 dissociation in hippocampal subfield coding as a result of memory reactivation.
836 Furthermore, we show that neural codes formed by hippocampal subregions not only
837 support simple recognition (Dimsdale-Zucker et al., 2018) or spatial memory (Leutgeb et

838 al., 2004, 2007), but also inference about the relationships among memories (see also
839 Schlichting et al., 2014). Inference decisions were faster and more accurate with
840 increasing similarity among indirectly items after learning in subiculum, indicating how
841 overlapping codes promote knowledge extraction beyond direct experience.

842 Our finding that subiculum representations track inference decisions may reflect
843 that subiculum is the output structure of the hippocampal circuit (O'Mara et al., 2001),
844 which plays a key role in recollection (Viskontas et al., 2009; Lindberg et al., 2017).
845 While subiculum showed evidence of learning-related differentiation for overlapping
846 pairs overall, our modeling data indicate that representational change in subiculum
847 reflects a continuum of responses. Increased integration (which can also be thought of
848 as less differentiation) promoted faster and more accurate inference. Our results
849 suggest that when memories are more integrated (or less differentiated), inference is
850 facilitated by retrieving a stored connection between indirectly related items (Shohamy
851 and Wagner, 2008; Schlichting et al., 2014); in contrast, differentiation might slow
852 inference between two separate traces that would need to be retrieved and recombined
853 at test (Koster et al., 2018).

854 Like subiculum, DG/CA_{2,3} exhibited learning-related differentiation of indirectly
855 related memory elements when memory reactivation was stronger during encoding.
856 However, it should be noted that DG/CA_{2,3} differentiation of overlapping memory
857 elements was only observed relative to the unrelated, across-triad baseline; there was
858 no change in similarity from pre- to post-learning for the indirectly-related items on their
859 own (**Fig. 4D** inset). This finding is consistent with prior work showing hippocampal
860 differentiation for related relative to unrelated events after learning (Favila et al., 2016;

861 Dimsdale-Zucker et al., 2018), while also controlling for baseline changes in similarity
862 that occur over time. Moreover, DG/CA_{2,3} showed evidence for memory integration
863 when memory reactivation was weaker during learning, suggesting the potential for
864 more nuanced representational dynamics in this region. For instance, memory
865 competition elicited by reactivation may have a non-monotonic relationship to
866 representational change in DG/CA_{2,3} (Ritvo et al., 2019). Stronger reactivation may
867 promote active differentiation; weaker or intermediate levels of reactivation may lead to
868 integration; and no reactivation may produce non-overlapping representations that are
869 separated via passive orthogonalization. This complex coding strategy could explain
870 why DG/CA_{2,3} shows evidence of differentiated (Kim et al., 2017) and integrated
871 (Schapiro et al., 2012) representations under different circumstances. Alternatively, our
872 results may reflect the use of a combined DG/CA_{2,3} region, the components of which
873 are thought to exhibit different transfer functions between environmental cues and
874 resulting memory representations (Yassa and Stark, 2011). The observed pattern of
875 results indicates that quantifying memory reactivation along with representational
876 change is necessary to fully understand how memory competition impacts
877 representation learning in DG/CA_{2,3}.

878 In summary, our empirical findings support a recently proposed computational
879 model of the hippocampal circuit (Schapiro et al., 2017); simulations from this model
880 suggest that CA₁ may represent relationships across events, whereas DG and CA₃
881 representations may emphasize differences between similar episodes. Our findings
882 align with these computational predictions, with CA₁ forming integrated representations
883 for similar memories, while DG/CA_{2,3} and subiculum differentiate those same

884 experiences. Additionally, we show that hippocampal representations support novel
885 inference, facilitating the discovery of unobserved relationships between distinct but
886 related experiences. The present work further shows that hippocampal subfield
887 dissociations are not a simple function of sensory input, but result from memory-based
888 competition during learning. Taken together, the present study advances our
889 understanding of how prior knowledge shapes how new events are represented within
890 the hippocampal circuit, providing an empirical test of key predictions of computational
891 models of hippocampal memory function.

892 **References**

- 893 Annis J, Miller BJ, Palmeri TJ (2017) Bayesian inference with Stan: a tutorial on adding
894 custom distributions. *Behav Res Methods* 49:863–886.
- 895 Avants BB, Tustison NJ, Song G, Cook PA, Klein A, Gee JC (2011) A reproducible
896 evaluation of ANTs similarity metric performance in brain image registration.
897 *Neuroimage* 54:2033–2044.
- 898 Bakker A, Kirwan CB, Miller M, Stark CEL (2008) Pattern separation in the human
899 hippocampal ca3 and dentate gyrus. *Science* 319:1640–1642.
- 900 Brainard DH (1997) The Psychophysics Toolbox. *Spat Vis* 10:433–436.
- 901 Brown SD, Heathcote A (2008) The simplest complete model of choice response time:
902 linear ballistic accumulation. *Cogn Psychol* 57:153–178.
- 903 Cox RW (1996) AFNI: software for analysis and visualization of functional magnetic
904 resonance neuroimages. *Comput Biomed Res* 29:162–173.
- 905 Dimsdale-Zucker HR, Ritchey M, Ekstrom AD, Yonelinas AP, Ranganath C (2018) CA₁
906 and CA₃ differentially support spontaneous retrieval of episodic contexts within
907 human hippocampal subfields. *Nat Commun* 9.
- 908 Duvernoy HM (1998) *The Human Hippocampus Functional Anatomy, Vascularization*
909 *and Serial Sections with MRI*. New York: Springer.
- 910 Eichenbaum H, Dudchenko P, Wood E, Shapiro M, Tanila H (1999) The hippocampus,
911 memory, and place cells: Is it spatial memory or a memory space? *Neuron* 23:209–
912 226.
- 913 Favila SE, Chanales AJH, Kuhl BA (2016) Experience-dependent hippocampal pattern
914 differentiation prevents interference during subsequent learning. *Nat Commun* 6:1–

- 915 10.
- 916 Guzowski JF, Knierim JJ, Moser EI (2004) Ensemble Dynamics of Hippocampal
917 Regions CA₃ and CA₁. *Neuron* 44:581–584.
- 918 Hanke M, Halchenko YO, Sederberg PB, Hanson SJ, Haxby J V., Pollmann S (2009)
919 PyMVPA: A python toolbox for multivariate pattern analysis of fMRI data.
920 *Neuroinformatics* 7:37–53.
- 921 Hsu NS, Schlichting ML, Thompson-Schill SL (2014) Feature diagnosticity affects
922 representations of novel and familiar objects. *J Cogn Neurosci* 26:2735–2749.
- 923 Hulbert JC, Norman KA (2015) Neural differentiation tracks improved recall of
924 competing memories following interleaved study and retrieval practice. *Cereb*
925 *Cortex* 25:3994–4008.
- 926 Jenkinson M (2003) Fast, automated, N-dimensional phase-unwrapping algorithm.
927 *Magn Reson Med* 49:193–197.
- 928 Kim G, Norman KA, Turk-Browne NB (2017) Neural differentiation of incorrectly
929 predicted memories. *J Neurosci* 37:2022–203.
- 930 Kirwan CB, Stark CEL (2007) Overcoming interference: An fMRI investigation of pattern
931 separation in the medial temporal lobe. *Learn Mem* 14:625–633.
- 932 Kleiner M, Brainard D, Pelli D, Ingling A, Murray R, Broussard C (2007) What’s new in
933 Psychtoolbox-3. *Perception* 36:1.
- 934 Koster R, Chadwick MJ, Chen Y, Berron D, Banino A, Düzel E, Hassabis D, Kumaran D
935 (2018) Big-loop recurrence within the hippocampal system supports integration of
936 information across episodes. *Neuron* 99:1342-1354.e6.
- 937 Kriegeskorte N, Goebel R, Bandettini P (2006) Information-based functional brain

- 938 mapping. *Proc Natl Acad Sci* 103:3863–3868.
- 939 Kriegeskorte N, Mur M, Bandettini P (2008) Representational similarity analysis -
940 connecting the branches of systems neuroscience. *Front Syst Neurosci* 2:4.
- 941 Kuhl BA, Rissman J, Chun MM, Wagner AD (2011) Fidelity of neural reactivation
942 reveals competition between memories. *Proc Natl Acad Sci* 108:5903–5908.
- 943 Lacy JW, Yassa M a, Stark SM, Muftuler LT, Stark CEL (2011) Distinct pattern
944 separation related transfer functions in human CA₃/dentate and CA₁ revealed using
945 high-resolution fMRI and variable mnemonic similarity. *Learn Mem* 18:15–18.
- 946 Lee I, Rao G, Knierim JJ (2004) A double dissociation between hippocampal subfields:
947 differential time course of CA₃ and CA₁ place cells for processing changed
948 environments. *Neuron* 42:803–815.
- 949 Leutgeb JK, Leutgeb S, Moser M-B, Moser EI (2007) Pattern separation in the dentate
950 gyrus and CA3 of the hippocampus. *Science* 315:961–966.
- 951 Leutgeb S, Leutgeb JK, Treves A, Moser M-B, Moser EI (2004) Distinct ensemble
952 codes in hippocampal areas CA₃ and CA₁. *Science* 305:1295–1298.
- 953 Lindberg O, Mårtensson G, Stomrud E, Palmqvist S, Wahlund LO, Westman E,
954 Hansson O (2017) Atrophy of the posterior subiculum is associated with memory
955 impairment, Tau- and A β pathology in non-demented individuals. *Front Aging*
956 *Neurosci* 9:1–12.
- 957 Mai J, Paxinos G, Voss T (2007) *Atlas of the Human Brain, Third*. Academic Press.
- 958 Marr D (1971) Simple memory: A theory for archicortex. *Philosophical Trans R Soc*
959 *London Ser B, Biol Sci* 262:23–81.
- 960 Morton NW, Schlichting ML, Preston AR. In press. Representations of common event

- 961 structure in medial temporal lobe and frontoparietal cortex support efficient
962 inference. *Proc Natl Acad Sci*.
- 963 Mumford JA, Turner BO, Ashby FG, Poldrack RA (2012) Deconvolving BOLD activation
964 in event-related designs for multivoxel pattern classification analyses. *Neuroimage*
965 59:2636–2643.
- 966 O'Mara SM, Commins S, Anderson M, Gigg J (2001) The subiculum: a review of form,
967 physiology and function. *Prog Neurobiol* 64:129–155.
- 968 Pelli DG (1997) The VideoToolbox software for visual psychophysics: transforming
969 numbers into movies. *Spat Vis* 10:437–442.
- 970 Polyn SM, Natu VS, Cohen JD, Norman KA (2005) Category-specific cortical activity
971 precedes retrieval during memory search. *Science* 310:1963–1966.
- 972 Potvin O, Doré FY, Goulet S (2009) Lesions of the dorsal subiculum and the dorsal
973 hippocampus impaired pattern separation in a task using distinct and overlapping
974 visual stimuli. *Neurobiol Learn Mem* 91:287–297.
- 975 Power JD, Barnes KA, Snyder AZ, Schlaggar BL, Petersen SE (2012) Spurious but
976 systematic correlations in functional connectivity MRI networks arise from subject
977 motion. *Neuroimage* 59:2142–2154.
- 978 Richter FR, Chanales AJH, Kuhl BA (2016) Predicting the integration of overlapping
979 memories by decoding mnemonic processing states during learning. *Neuroimage*
980 124:323–335.
- 981 Ritvo VJH, Turk-Browne NB, Norman KA (2019) Nonmonotonic plasticity: How memory
982 retrieval drives learning. *Trends Cogn Sci* 23:726–742.
- 983 Rumelhart DE, Hinton GE, Williams RJ (1986) Learning internal representations by

- 984 error propagation. In: *Parallel Distributed Processing: Explorations in the*
985 *Microstructure of Cognition (Foundations, Vol. 1)*, pp 318–362. MIT Press.
- 986 Schapiro AC, Kustner L V., Turk-Browne NB (2012) Shaping of object representations
987 in the human medial temporal lobe based on temporal regularities. *Curr Biol*
988 22:1622–1627.
- 989 Schapiro AC, Turk-Browne NB, Botvinick MM, Norman KA (2017) Complementary
990 learning systems within the hippocampus: A neural network modeling approach to
991 reconciling episodic memory with statistical learning. *Philosophical Trans R Soc B*.
- 992 Schlichting ML, Mack ML, Guarino KF, Preston AR (2019) Performance of semi-
993 automated hippocampal subfield segmentation methods across ages in a pediatric
994 sample. *Neuroimage* 191:49–67.
- 995 Schlichting ML, Mumford JA, Preston AR (2015) Learning-related representational
996 changes reveal dissociable integration and separation signatures in the
997 hippocampus and prefrontal cortex. *Nat Commun* 6:8151.
- 998 Schlichting ML, Preston AR (2014) Memory reactivation during rest supports upcoming
999 learning of related content. *Proc Natl Acad Sci* 111:15845–15850.
- 1000 Schlichting ML, Preston AR (2015) Memory integration: Neural mechanisms and
1001 implications for behavior. *Curr Opin Behav Sci* 1:1–8.
- 1002 Schlichting ML, Zeithamova D, Preston AR (2014) CA₁ subfield contributions to memory
1003 integration and inference. *Hippocampus* 24:1248–1260.
- 1004 Shohamy D, Wagner AD (2008) Integrating memories in the human brain:
1005 Hippocampal-midbrain encoding of overlapping events. *Neuron* 60:378–389.
- 1006 Vazdarjanova A, Guzowski JF (2004) Differences in hippocampal neuronal population

- 1007 responses to modifications of an environmental context: Evidence for distinct, yet
1008 complementary, functions of CA₃ and CA₁ ensembles. *J Neurosci* 24:6489–6496.
- 1009 Vehtari A, Gelman A, Gabry J (2017) Practical Bayesian model evaluation using leave-
1010 one-out cross-validation and WAIC. *Stat Comput* 27:1413–1432.
- 1011 Vieweg P, Stangl M, Howard LR, Wolbers T (2015) Changes in pattern completion - A
1012 key mechanism to explain age-related recognition memory deficits? *Cortex* 64:343–
1013 351.
- 1014 Viskontas I V., Carr VA, Engel SA, Knowlton BJ (2009) The neural correlates of
1015 recollection: hippocampal activation declines as episodic memory fades.
1016 *Hippocampus* 19:265–272.
- 1017 West MJ, Gundersen HJG (1990) Unbiased stereological estimation of the number of
1018 neurons in the human hippocampus. *J Comp Neurol* 296:1–22.
- 1019 Winkler AM, Ridgway GR, Webster MA, Smith SM, Nichols TE (2014) Permutation
1020 inference for the general linear model. *Neuroimage* 92:381–397.
- 1021 Yassa MA, Stark CEL (2011) Pattern separation in the hippocampus. *Trends Neurosci*
1022 34:515–525.
- 1023 Zeithamova D, Dominick AL, Preston AR (2012) Hippocampal and ventral medial
1024 prefrontal activation during retrieval-mediated learning supports novel inference.
1025 *Neuron* 75:168–179.
- 1026 Zeithamova D, Gelman BD, Frank L, Preston AR (2018) Abstract representation of
1027 prospective reward in the hippocampus. *J Neurosci* 38:10093–10101.
- 1028 Zhang Y, Brady M, Smith S (2001) Segmentation of brain MR images through a hidden
1029 Markov random field model and the expectation-maximization algorithm. *IEEE*

1030 Trans Med Imaging 20:45–57.

1031

1032 **Figure Captions**

1033 **Figure 1.** Experimental design. **(A)** Schematic of the behavioral task. Participants were
1034 first exposed to individually presented pictures (faces, scenes, and novel objects) that
1035 would later become indirectly related through associative learning (A and C items).
1036 Then, participants learned to associate initial pairs (face-shape or scene-shape AB
1037 associations) and were scanned while learning overlapping pairs (shape-object BC
1038 associations). Participants were scanned again in a post-exposure phase while they
1039 viewed the same items from pre-exposure (A and C items). Participants then completed
1040 an across-episode inference task. Finally, participants completed a localizer task in
1041 which they viewed individually presented faces, scenes, objects, and shapes in a
1042 blocked design. **(B)** Visual similarity manipulation. The similarity of the shared B item
1043 across pairs was parametrically manipulated. In this example, the top shape would have
1044 been seen in the initial AB pairs, while the bottom row depicts the different shape
1045 morphs that could be seen when learning the overlapping BC pairs. The linking B item
1046 presented during overlapping pair learning could either be an exact match to the B item
1047 presented during initial (AB) pair learning, a high similarity or low similarity morph, or
1048 new (i.e., non-overlapping) item. **(C)** Subjective similarity of shape stimuli used for B
1049 linking items. An independent sample of participants rated visual similarity between
1050 parent shapes and shape morphs presented side by side using a 5-point Likert scale (1
1051 = not at all similar, 5 = very similar). Significance of paired *t*-tests are shown with
1052 asterisks (*) for $p < 0.05$. Error bars represent \pm standard error of the mean.

1053

1054 **Figure 2.** Behavioral performance. **(A)** Overlapping pair (BC) test accuracy and **(B)**
1055 response time (correct trials only) by learning block for each similarity condition. **(C)**
1056 Across-episode (AC) inference accuracy and **(D)** response time (correct trials only) for
1057 each similarity condition. Significance of paired *t*-tests are shown with asterisks (*) for p
1058 < 0.05 . Error bars represent \pm standard error of the mean. Dotted lines indicate chance
1059 performance on the 3 ACF tests.

1060

1061 **Figure 3.** Memory reactivation during overlapping pair learning. **(A)** Results of the
1062 searchlight analysis identifying regions where classifier evidence for A item reactivation
1063 exceeded baseline (i.e., evidence during new, non-overlapping pairs) when participants
1064 were learning the overlapping BC pairs. **(B)** Evidence for reactivation of A items as a
1065 function of stimulus category (face and scene) during overlapping pair learning for each
1066 of the regions identified in **(A)**. Error bars represent \pm standard error of the mean. **(C)**
1067 Results of the searchlight analysis identifying regions where classifier evidence for A
1068 item reactivation varied with visual similarity of the linking B item (exact match $>$ high
1069 and low similarity). One cluster in left parietal cortex overlapped with the cluster
1070 identified in the searchlight analysis comparing reactivation to baseline (**A**, leftmost
1071 image); the other cluster extended into occipital cortex. All searchlight clusters are
1072 displayed on the 1mm MNI 152 anatomical template.

1073

1074 **Figure 4.** Assessing learning-related representational change as a function of memory
1075 reactivation during learning. **(A)** Predictions for memory formation through associative
1076 learning. Prior to learning, individual A and C items in the pre-exposure phase do not

1077 share any relationships. After learning, the representations of A and C items may shift
1078 as a function of their shared relationships with B items. We tested for two neural
1079 outcomes; in the case of differentiation, the neural patterns for indirectly related A and C
1080 items are predicted to be less similar in the post-exposure phase relative to the pre-
1081 learning representations. In contrast, for memory integration, the neural similarity of
1082 indirectly related A and C items are predicted to increase from pre- to post-learning,
1083 reflecting the formation of overlapping neural codes linking elements experienced
1084 across events. **(B)** Four searchlight contrasts were used to determine whether memory
1085 representation in hippocampal subfields varied with memory reactivation strength during
1086 learning. Two of the searchlights identified regions in which differentiation or integration
1087 occurred across all degrees of reactivation strength. Another set of searchlights
1088 identified regions in which neural coding varied as a function of reactivation. **(C)**
1089 Learning-related representational change in hippocampus. Subregions of DG/CA_{2,3} and
1090 subiculum showed differentiation of the indirectly related elements of overlapping
1091 memories, but only when reactivation was stronger during learning. In contrast, a
1092 subregion of CA₁ showed evidence of memory integration, but again only when
1093 reactivation was stronger during overlapping pair learning. Hippocampal regions are
1094 depicted on an open source high-resolution group T2 template created for hippocampal
1095 subfield analyses (Schlichting et al., 2019). **(D)** Neural similarity change in the clusters
1096 identified in the searchlight analysis **(C)** after reverse-normalization to native space,
1097 confirming the predicted pattern of results from **(B)**. The inset displays the same data
1098 separately for the within-triad and across-triad similarity measures prior to calculating
1099 the difference scores. Note that because this analysis is based on voxels identified in

1100 the searchlight analysis, it is not fully independent. Error bars represent \pm standard error
1101 of the mean.

1102

1103 **Figure 5.** Results of the multilevel response time model used to examine relationships
1104 between neural measures and AC inference performance. **(A)** Fit of response time
1105 model to accuracy on individual AC inference trials. **(B)** Fit of response time model to
1106 trial-level inference response times. **(C)** We examined whether reactivation of related
1107 memories during BC study or neural similarity change (Δ) in hippocampal subfields after
1108 learning predicted trial-level variability in AC inference performance (i.e., the slope of
1109 the drift rate from the model). Negative values indicate a decrease in the neural
1110 measures predicted faster and more accurate inference, while positive values indicate
1111 an increase in the neural measures predicted better inference. Reactivation of related
1112 memories and representational change within subiculum predicted improved AC
1113 inference performance. Bars indicate 95% high-density intervals of posterior parameter
1114 estimates.









