

¹ A novel framework for linking dynamic experiences to
² memories reveals event-like structure in naturalistic
³ episodic recall

⁴ Andrew C. Heusser, Paxton C. Fitzpatrick, and Jeremy R. Manning

Department of Psychological and Brain Sciences

Dartmouth College, Hanover, NH 03755, USA

Corresponding author: jeremy.r.manning@dartmouth.edu

⁵ July 21, 2019

⁶ **Abstract**

⁷ Our life experiences unfold over time and the temporal dynamics of their contents form
⁸ unique *experience trajectories*. Within this geometric framework, one can compare the shape of the
⁹ trajectory formed by an experience to that defined by our later remembering of that experience.
¹⁰ We propose a framework for mapping naturalistic experiences onto geometric spaces that char-
¹¹ acterize how they unfold over time. We apply this approach to a naturalistic memory experiment
¹² which had participants view and recount a video. We find that the video and subsequent recall
¹³ share an event-like structure, and that the “shapes” of the trajectories formed by participants’
¹⁴ recounts were all highly similar to that of the original video. Further, participants’ that re-
¹⁵ counted the events with a high level of precision and in a distinctive way had better subsequent
¹⁶ memory performance. Lastly, we identified a network of brain structures that are sensitive to
¹⁷ the shapes of our ongoing experiences, and an overlapping network that is sensitive to how
¹⁸ we will later remember those experiences. These results highlight the rich event-like structure

19 of memories for our life experiences, and introduce a novel framework for mapping between
20 dynamic life experiences and their mnemonic counterparts.

21 **Introduction**

22 What does it mean to *remember* something? In traditional episodic memory experiments (e.g.,
23 list-learning or trial-based experiments; Murdock, 1962; Kahana, 1996), remembering is often cast
24 as a discrete and binary operation: each studied item may be separated from the rest of one's
25 experiences, and that item may be labeled as having been recalled versus forgotten. More nuanced
26 studies might incorporate self-reported confidence measures as a proxy for memory strength, or
27 ask participants to discriminate between "recollecting" the (contextual) details of an experience
28 or having a general feeling of "familiarity" (Yonelinas, 2002). Using well-controlled trial-based
29 experimental designs, the field has amassed a wealth of valuable information regarding human
30 episodic memory. However, there are fundamental properties of our memories that trial-based
31 experiments are not well suited to capture (for review also see Koriat and Goldsmith, 1994; Huk
32 et al., 2018). First, our memories are continuous, rather than discrete: removing a (naturalistic)
33 event from the context in which it occurs can substantially change its meaning. Second, the specific
34 language used to describe an experience have little bearing on whether the experience should be
35 considered to have been "remembered." Asking whether the rememberer has precisely reproduced
36 a specific set of words to describe a given experience is nearly orthogonal to whether they were
37 actually able to remember it. In classic (e.g., list-learning) memory studies, by contrast, counting
38 the number or proportion of precise recalls is often a primary metric of assessing the quality of
39 participants' memories. Third, one might remember the *essence* (or shape) of an experience but
40 forget (or neglect to recount) particular details. Capturing the essence of what happened is typically
41 the main "point" of recounting a memory to a listener.

42 How might one go about formally characterizing the shape of an experience, or whether that
43 shape has been recovered by the rememberer? Any given moment of an experience derives
44 meaning from surrounding moments, as well as from longer-range temporal associations (e.g.,

45 Lerner et al., 2011). Therefore, the timecourse describing how an event unfolds is fundamental
46 to its overall meaning. Further, this hierarchy formed by our subjective experiences at different
47 timescales defines a *context* for each new moment (e.g., Howard and Kahana, 2002; Howard et al.,
48 2014), and plays an important role in how we interpret that moment and remember it later (for
49 review see Manning et al., 2015). Our memory systems can then leverage these associations to
50 form predictions that help guide our behaviors (Ranganath and Ritchey, 2012). For example,
51 as we navigate the world, the features of our subjective experiences tend to change gradually
52 (e.g., the room or situation we are in is strongly temporally autocorrelated), allowing us to form
53 stable estimates of our current situation and behave accordingly (Zacks et al., 2007; Zwaan and
54 Radvansky, 1998).

55 Although our experiences most often change gradually, they also occasionally change sud-
56 denly (e.g., when we walk through a doorway; Radvansky and Zacks, 2017). Prior research
57 suggests that these sharp transitions (termed *event boundaries*) during an experience help to dis-
58 cretize our experiences into *events* (Radvansky and Zacks, 2017; Brunec et al., 2018; Heusser et al.,
59 2018a; Clewett and Davachi, 2017; Ezzyat and Davachi, 2011; DuBrow and Davachi, 2013). The
60 interplay between the stable (within event) and transient (across event) temporal dynamics of an
61 experience also provides a potential framework for transforming experiences into memories that
62 distill those experiences down to their essence. For example, prior work has shown that event
63 boundaries can influence how we learn sequences of items (Heusser et al., 2018a; DuBrow and
64 Davachi, 2013), navigate (Brunec et al., 2018), and remember and understand narratives (Zwaan
65 and Radvansky, 1998; Ezzyat and Davachi, 2011). Prior research has implicated the hippocampus
66 and the medial prefrontal cortex as playing a critical role in transforming experiences into stuctured
67 and consolidated memories (Tompry and Davachi, 2017).

68 Here we sought to examine how the temporal dynamics of a “naturalistic” experience were
69 reflected in participants’ later memories of that experience. We analyzed an open dataset that
70 comprised behavioral and functional Magnetic Resonance Imaging (fMRI) data collected as par-
71 ticipants viewed and then verbally recalled an episode of the BBC television series *Sherlock* (Chen
72 et al., 2017). We developed a computational framework for characterizing the temporal dynamics

73 of the moment-by-moment content of the episode and of participants' verbal recalls. Specifically,
74 we use topic modeling (Blei et al., 2003) to characterize the thematic conceptual (semantic) content
75 present in each moment of the episode and recalls, and we use Hidden Markov Models (Rabiner,
76 1989; Baldassano et al., 2017) to discretize the evolving semantic content into events. In this way,
77 we cast naturalistic experiences (and recalls of those experiences) as *trajectories* that describe how
78 the experiences evolve over time. Under this framework, successful remembering entails verbally
79 "traversing" the topic trajectory of the original episode, thereby reproducing the shape of the
80 original experience. Comparing the shapes of the topic trajectories of the original episode and of
81 participants' retellings of the episode reveals which aspects of the episode were preserved (or lost)
82 in the translation into memory.

83 Here, we introduce a novel content-based framework for linking dynamic experiences to mem-
84 ories. We hypothesized that the "shape" and the event-like structure of an experience would
85 be preserved in memory. Further, we tested 1) whether the *precision* with which the participant
86 recounts each event and 2) how *distinctive* each recall event is (relative to the other recalled events)
87 relates to overall memory performance. Last, we identify networks of brain structures whose
88 responses (as participants watched the episode) reflected the shape of the episode, and how par-
89 ticipants would later recount the episode.

90 Results

91 To characterize the shape of the *Sherlock* episode participants watched and their subsequent re-
92 countings of the episode, we used a topic model (Blei et al., 2003) to discover the latent thematic
93 content in the video. Topic models take as inputs a vocabulary of words to consider and a collection
94 of text documents; they return as output two matrices. The first output is a *topics matrix* whose
95 rows are topics (latent themes) and whose columns correspond to words in the vocabulary. The
96 entries of the topics matrix define how each word in the vocabulary is weighted by each discovered
97 topic. For example, a detective-themed topic might weight heavily on words like "crime," and
98 "search." The second output is a *topic proportions matrix*, with one row per document and one

99 column per topic. The topic proportions matrix describes which mixture of topics is reflected in
100 each document.

101 Chen et al. (2017) collected hand-annotated information about each of 1000 (manually identified)
102 scenes spanning the roughly 45 minute video used in their experiment. This information included:
103 a brief narrative description of what was happening; whether the scene took place indoors vs.
104 outdoors; names of any characters on the screen; names of any characters who were in focus in
105 the camera shot; names of characters who were speaking; the location where the scene took place;
106 the camera angle (close up, medium, long, etc.); whether or not background music was present;
107 and other similar details (for a full list of annotated features see *Methods*). We took from these
108 annotations the union of all unique words (excluding stop words, such as “and,” “or,” “but,” etc.)
109 across all features and scenes as the “vocabulary” for the topic model. We then concatenated the
110 sets of words across all features contained in overlapping 50-scene sliding windows, and treated
111 each 50-scene sequence as a single “document” for the purposes of fitting the topic model. Next,
112 we fit a topic model with (up to) $K = 100$ topics to this collection of documents. We found that 27
113 unique topics (with non-zero weights) were sufficient to describe the time-varying content of the
114 video (see *Methods*; Figs. 1, S2). Note that our approach is similar in some respects to Dynamic Topic
115 Models (Blei and Lafferty, 2006), in that we sought to characterize how the thematic content of the
116 episode evolved over time. However, whereas Dynamic Topic Models are designed to characterize
117 how the properties of *collections* of documents change over time, our sliding window approach
118 allows us to examine the topic dynamics within a single document (or video). Specifically, our
119 approach yielded (via the topic proportions matrix) a single *topic vector* for each timepoint of the
120 episode (we set timepoints to match the acquisition times of the 1976 fMRI volumes collected as
121 participants viewed the episode).

122 The topics we found were heavily character-focused (e.g., the top-weighted word in each topic
123 was nearly always a character) and could be roughly divided into themes that were primarily
124 Sherlock Holmes-focused (Sherlock is the titular character); primarily John Watson-focused (John
125 is Sherlock’s close confidant and assistant); or that involved Sherlock and John interacting (Fig. S2).
126 Several of the topics were highly similar, which we hypothesized might allow us to distinguish



Figure 1: Methods overview. We used hand-annotated descriptions of each moment of video to fit a topic model. Three example video frames and their associated descriptions are displayed (top two rows). Participants later recalled the video (in the third row, we show example recalls of the same three scenes from participant 13). We used the topic model (fit to the annotations) to estimate topic vectors for each moment of video and each sentence the participants recalled. Example topic vectors are displayed in the bottom row (blue: video annotations; green: example participant’s recalls). Three topic dimensions are shown (the highest-weighted topics for each of the three example scenes, respectively). We also show the ten highest-weighted words for each topic. Figure S2 provides a full list of the top 10 words from each of the discovered topics.

127 between subtle narrative differences (if the distinctions between those overlapping topics were
128 meaningful; also see Fig. S3). The topic vectors for each timepoint were *sparse*, in that only a small
129 number (usually one or two) of topics tended to be “active” in any given timepoint (Fig. 2A).
130 Further, the dynamics of the topic activations appeared to exhibit *persistance* (i.e., given that a
131 topic was active in one timepoint, it was likely to be active in the following timepoint) along with
132 *occasional rapid changes* (i.e., occasionally topics would appear to spring into or out of existence).
133 These two properties of the topic dynamics may be seen in the block diagonal structure of the
134 timepoint-by-timepoint correlation matrix (Fig. 2B). Following Baldassano et al. (2017), we used a
135 Hidden Markov Model (HMM) to identify the *event boundaries* where the topic activations changed
136 rapidly (i.e., at the boundaries of the blocks in the correlation matrix; event boundaries identified
137 by the HMM are outlined in yellow). Part of our model fitting procedure required selecting an
138 appropriate number of “events” to segment the timeseries into. We used an optimization procedure
139 to identify the number of events that maximized within-event stability while also minimizing
140 across-event correlations (see *Methods* for additional details). To create a stable “summary” of the
141 video, we computed the average topic vector within each event (Fig. 2C).

142 Given that the time-varying content of the video could be segmented cleanly into discrete
143 events, we wondered whether participants’ recalls of the video also displayed a similar structure.
144 We applied the same topic model (already trained on the video annotations) to each participant’s
145 recalls. Analogous to how we analyzed the time-varying content of the video, to obtain similar
146 estimates for participants’ recalls, we treated each (overlapping) 10 sentence “window” of their
147 transcript as a “document” and then computed the most probable mix of topics reflected in each
148 timepoint’s sentences. This yielded, for each participant, a number-of-sentences by number-of-
149 topics topic proportions matrix that characterized how the topics identified in the original video
150 were reflected in the participant’s recalls. Note that an important feature of our approach is
151 that it allows us to compare participant’s recalls to events from the original video, despite that
152 different participants may have used different language to describe the same event, and that those
153 descriptions may not match the original annotations. This is a huge benefit of projecting the video
154 and recalls into a shared “topic” space. An example topic proportions matrix from one participant’s

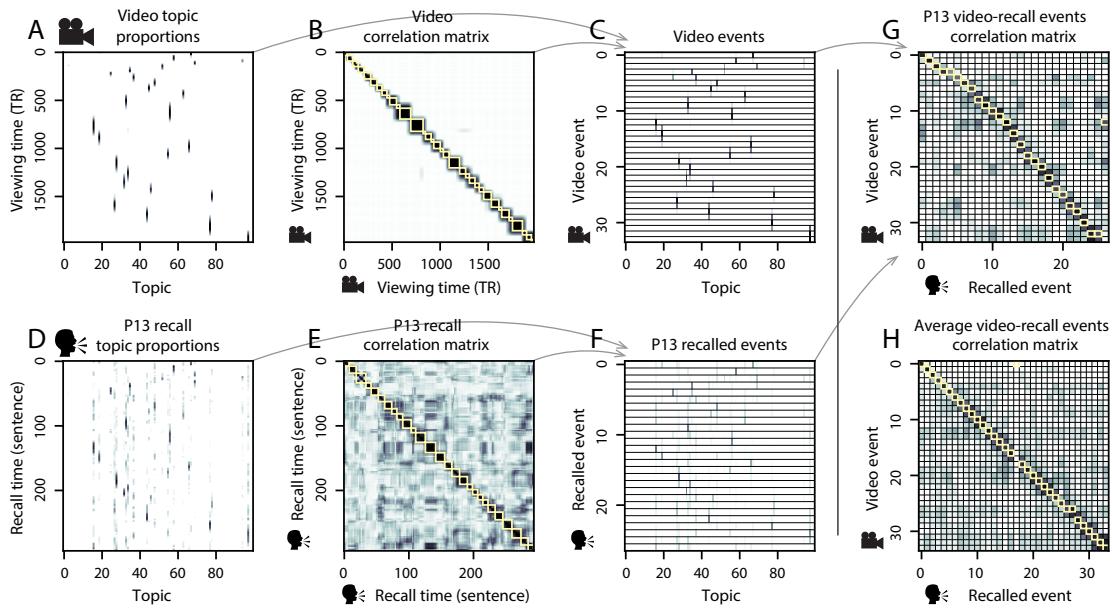


Figure 2: Modelling naturalistic stimuli and recalls. All panels: darker colors indicate greater values; range: [0, 1]. **A.** Topic vectors ($K = 100$) for each of the 1976 video timepoints. **B.** Timepoint-by-timepoint correlation matrix of the topic vectors displayed in Panel A. Event boundaries detected by the HMM are denoted in yellow (34 events detected). **C.** Average topic vectors for each of the 34 video events. **D.** Topic vectors for each of 294 sentences spoken by an example participant while recalling the video. **E.** Timepoint-by-timepoint correlation matrix of the topic vectors displayed in Panel D. Event boundaries detected by the HMM are denoted in yellow (27 events detected). **F.** Average topic vectors for each of the 27 recalled events from the example participant. **G.** Correlations between the topic vectors for every pair of video events (Panel C) and recalled events (from the example participant; Panel F). For similar plots for all participants see Figure S5. **H.** Average correlations between each pair of video events and recalled events (across all 17 participants). To create the figure, each recalled event was assigned to the video event with the most correlated topic vector (yellow boxes in panels G and H). The heat maps in each panel were created using Seaborn (Waskom et al., 2016).

155 recalls is shown in Figure 2D.

156 Although the example participant's recall topic proportions matrix has some visual similarity to
157 the video topic proportions matrix, the time-varying topic proportions for the example participant's
158 recalls are not as sparse as for the video (e.g., compare Figs. 2A and D). Similarly, although there
159 do appear to be periods of stability in the recall topic dynamics (e.g., most topics are active or
160 inactive over contiguous blocks of time), the overall timecourses are not as cleanly delineated as
161 the video topics are. To examine these patterns in detail, we computed the timepoint-by-timepoint
162 correlation matrix for the example participant's recall topic proportions (Fig. 2E). As in the video
163 correlation matrix (Fig. 2B), the example participant's recall correlation matrix has a strong block
164 diagonal structure, indicating that their recalls are discretized into separated events. As for the
165 video correlation matrix, we can use an HMM, along with the aforementioned number-of-events
166 optimization procedure (also see *Methods*) to determine how many events are reflected in the
167 participant's recalls and where specifically the event boundaries fall (outlined in yellow). We
168 carried out a similar analysis on all 17 participants' recall topic proportions matrices (Fig. S4).

169 Two clear patterns emerged from this set of analyses. First, although every individual partic-
170 ipant's recalls could be segmented into discrete events (i.e., every individual participant's recall
171 correlation matrix exhibited clear block diagonal structure; Fig. S4), each participant appeared to
172 have a unique *recall resolution*, reflected in the sizes of those blocks. For example, some participants'
173 recall topic proportions segmented into just a few events (e.g., Participants P1, P4, and P15), while
174 others' recalls segmented into many shorter duration events (e.g., Participants P12, P13, and P17).
175 This suggests that different participants may be recalling the video with different levels of detail-
176 e.g., some might touch on just the major plot points, whereas others might attempt to recall every
177 minor scene. The second clear pattern present in every individual participant's recall correlation
178 matrix is that, unlike in the video correlation matrix, there are substantial off-diagonal correlations
179 in participant's recalls. Whereas each event in the original video (was largely) separable from the
180 others (Fig. 2B), in transforming those separable events into memory participants appear to be
181 integrating *across* different events, blending elements of previously recalled and not-yet-recalled
182 events into each newly recalled event (Figs. 2D, S4; also see Manning et al., 2011; Howard et al.,

183 2012).

184 The above results indicate that both the structure of the original video and participants' recalls
185 of the video exhibit event boundaries that can be identified automatically by characterizing the
186 dynamic content using a shared topic model and segmenting the content into events using HMMs.
187 Next we asked whether some correspondence might be made between the specific content of
188 the events the participants experienced in the video, and the events they later recalled. One
189 approach to linking the experienced (video) and recalled events is to label each recalled event as
190 matching the video event with the most similar (i.e., most highly correlated) topic vector (Figs. 2G,
191 S5). This yields a sequence of "presented" events from the original video, and a sequence of
192 (potentially differently ordered) "recalled" events for each participant. Analogous to classic list-
193 learning studies, we can then examine participants' recall sequences by asking which events
194 they tended to recall first (probability of first recall; Fig. 3A; Welch and Burnett, 1924; Postman
195 and Phillips, 1965; Atkinson and Shiffrin, 1968); how participants most often transition between
196 recalls of the events as a function of the temporal distance between them (lag-conditional response
197 probability; Fig. 3B; Kahana, 1996); and which events they were likely to remember overall (serial
198 position recall analyses; Fig. 3C; Murdock, 1962). Interestingly, for 2 of the analyses (probability
199 of first recall and lag-conditional response probability curves) we observe patterns comparable to
200 classic effects from the list-learning literature. Namely, a higher probability of initiating recall with
201 the first event in the sequence (Fig. 3A) and a higher probability of transitioning to neighboring
202 events with a forward asymmetric bias (Fig. 3C). In contrast, we do not observe a pattern comparable
203 to the serial position effect (Fig. 3C), but rather we see higher memory for specific events distributed
204 somewhat evenly throughout the video.

205 Statistical models of memory studies often treat memory recalls as binary (e.g. the item was
206 recalled or not) and independent events. However, our framework produces a content-based model
207 of individual stimulus and recall events, allowing for direct quantitative comparison between all
208 stimulus and recall events, as well as between the recall events themselves. Leveraging these
209 content-based models of the stimulus/recall events, we developed 2 novel metrics for quantifying
210 naturalistic memory representations: precision and distinctiveness. We define precision as the

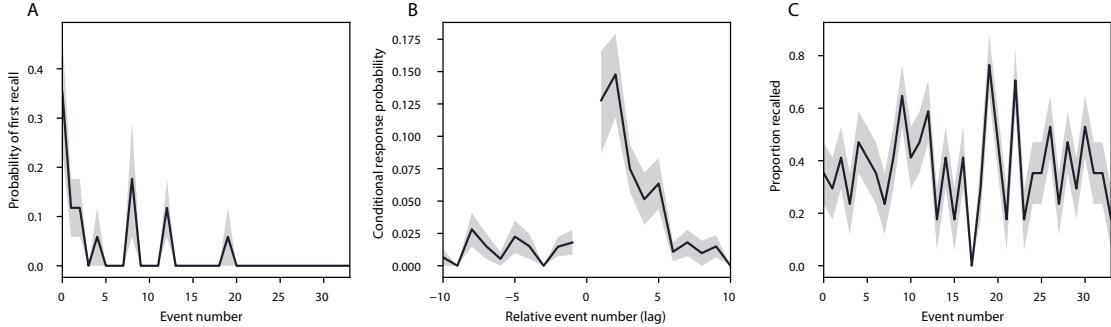


Figure 3: Naturalistic extensions of classic list-learning memory analyses. A. The probability of first recall as a function of the serial position of the event in the video. B. The probability of recalling each event, conditioned on having most recently recalled the event *lag* events away in the video. C. The proportion of participants who recalled each event, as a function of the serial position of the events in the video. All panels: error bars denote bootstrap-estimated standard error of the mean.

211 average correlation between each recall event and the maximally correlated video event (Fig. 4).
 212 Participants whose recall events are more veridical descriptions of what happened in the video
 213 event will presumably have higher precision scores. We find that across participants, a higher
 214 precision score is correlated to both hand annotated memory performance (Pearson's $r(15) =$
 215 $0.6, p = 0.011$) as well as the number of recall events estimated by our model (Pearson's $r(15) =$
 216 $0.64, p = 0.005$). A second novel metric we introduce here is distinctiveness, or how unique the
 217 recall description was to each video event. We define distinctiveness as 1 minus the average of
 218 all non-matching recall events from the video-recall correlation matrix. We hypothesized that
 219 participants who recounted events in a more distinctive way would display better overall memory.
 220 Similar to precision, we find that the more distinct participants recalls are (on average), the more
 221 they remembered (hand-annotated memory: Pearson's $r(15) = 0.83, p < 0.001$ and model derived
 222 memory: Pearson's $r(15) = 0.71, p = 0.001$). In summary, using two novel metrics afforded by our
 223 approach, we find that participants whose recalls are both more precise and distinct remember
 224 more content.

225 The prior analyses leverage the correspondence between the 100-dimensional topic proportion
 226 matrices for the video and participants' recalls to characterize recall. However, it is difficult
 227 to gain deep insights into that content solely by examining the topic proportion matrices (e.g.,

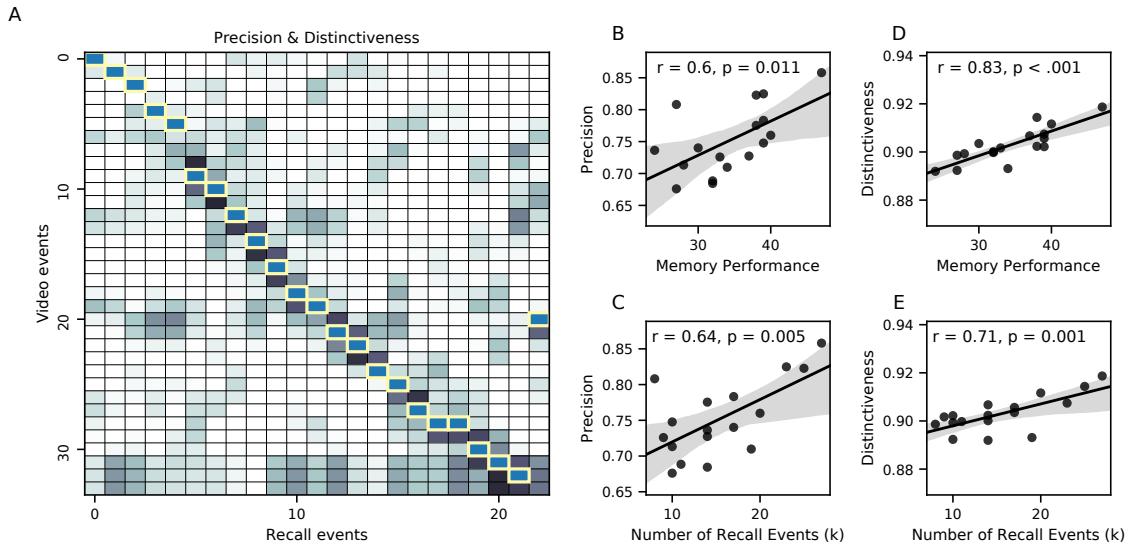


Figure 4: Novel content-based metrics of naturalistic memory: precision and distinctiveness. **A.** A video-recall correlation matrix for a representative participant (17). The yellow boxes highlight the maximum correlation in each column. Precision was computed as the average of the maximum correlation in each column. On the other hand, distinctiveness was defined as the average of everything except for the maximum correlation in each column. **B.** The (Pearson's) correlation between precision and hand-annotated memory performance. **C.** The correlation between precision and the number of events recovered by the model (k). **D.** The correlation between distinctiveness and hand-annotated memory performance. **E.** The correlation between distinctiveness and the number of events recovered by the model (k).

228 Figs. 2A, D) or the corresponding correlation matrices (Figs. 2B, E, S4). To visualize the time-
229 varying high-dimensional content in a more intuitive way (Heusser et al., 2018b) we projected the
230 topic proportions matrices onto a two-dimensional space using Uniform Manifold Approximation
231 and Projection (UMAP; McInnes and Healy, 2018). In this lower-dimensional space, each point
232 represents a single video or recall event, and the distances between the points reflect the distances
233 between the events' associated topic vectors (Fig. 5). In other words, events that are near to each
234 other in this space are more semantically similar.

235 Visual inspection of the video and recall topic trajectories reveals a striking pattern. First,
236 the topic trajectory of the video (which reflects its dynamic content; Fig. 5A) is captured nearly
237 perfectly by the averaged topic trajectories of participants' recalls (Fig. 5B). To assess the consistency
238 of these recall trajectories across participants, we asked: given that a participant's recall trajectory
239 had entered a particular location in topic space, could the position of their *next* recalled event
240 be predicted reliably? For each location in topic space, we computed the set of line segments
241 connecting successively recalled events (across all participants) that intersected that location (see
242 *Methods* for additional details). We then computed (for each location) the distribution of angles
243 formed by the lines defined by those line segments and a fixed reference line (the *x*-axis). Rayleigh
244 tests revealed the set of locations in topic space at which these across-participant distributions
245 exhibited reliable peaks (blue arrows in Fig. 5B reflect significant peaks at $p < 0.05$, corrected). We
246 observed that the locations traversed by nearly the entire video trajectory exhibited such peaks.
247 In other words, participants exhibited similar trajectories that also matched the trajectory of the
248 original video (Fig. 5C). This is especially notable when considering the fact that the number of
249 events participants recalled (dots in Fig. 5C) varied considerably across people, and that every
250 participant used different words to describe what they had remembered happening in the video.
251 Differences in the numbers of remembered events appear in participants' trajectories as differences
252 in the sampling resolution along the trajectory. We note that this framework also provides a
253 means of detangling classic "proportion recalled" measures (i.e., the proportion of video events
254 referenced in participants' recalls) from participants' abilities to recapitulate the full shape of the
255 original video (i.e., the similarity in the shape of the original video trajectory and that defined by

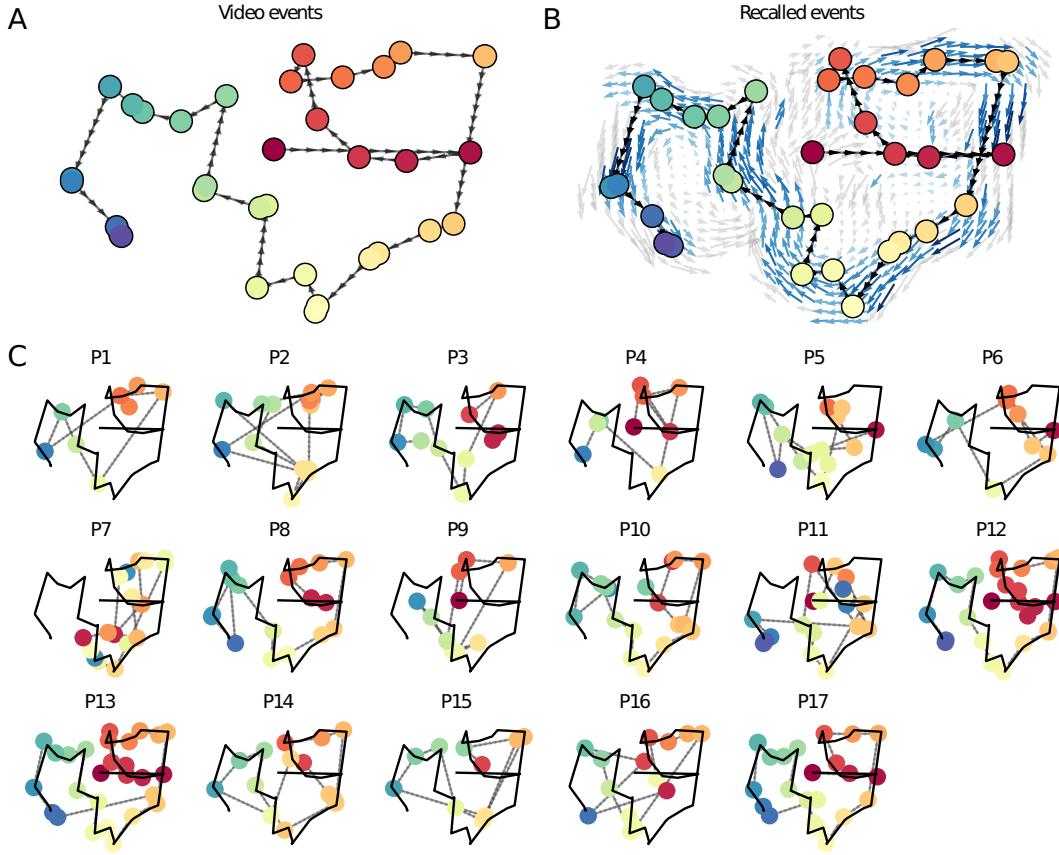


Figure 5: Trajectories through topic space capture the dynamic content of the video and recalls. All panels: the topic proportion matrices have been projected onto a shared two-dimensional space using UMAP. **A.** The two-dimensional topic trajectory taken by the episode of *Sherlock*. Each dot indicates an event identified using the HMM (see *Methods*); the dot colors denote the order of the events (early events are in red; later events are in blue), and the connecting lines indicate the transitions between successive events. **B.** The average two-dimensional trajectory captured by participants' recall sequences, with the same format and coloring as the trajectory in Panel A. To compute the event positions, we matched each recalled event with an event from the original video (see *Results*), and then we averaged the positions of all events with the same label. The arrows reflect the average transition direction through topic space taken by any participants whose trajectories crossed that part of topic space; blue denotes reliable agreement across participants via a Rayleigh test ($p < 0.05$, corrected). **C.** The recall topic trajectories (gray) taken by each individual participant (P1–P17). The video's trajectory is shown in black for reference. (Same format and coloring as Panel A.)

256 each participant's recounting of the video).

257 Because our analysis framework projects the dynamic video content and participants' recalls
258 onto a shared topic space, and because the dimensions of that space are known (i.e., each topic
259 dimension is a set of weights over words in the vocabulary; Fig. S2), we can examine the topic
260 trajectories to understand which specific content was remembered well (or poorly). For each video
261 event, we can ask: what was the average correlation (across participants) between the video event's
262 topic vector and the closest matching recall event topic vectors from each participant? This yields a
263 single correlation coefficient for each video event, describing how closely participants' recalls of the
264 event tended to reliably capture its content (Fig. 6A). (We also examined how different comparisons
265 between each video event's topic vector and the corresponding recall event topic vectors related
266 to hand-annotated characterizations of memory performance; see *Supporting Information*). Given
267 this summary of which events were recalled reliably (or not), we next asked whether the better-
268 remembered or worse-remembered events tended to reflect particular topics. We computed a
269 weighted average of the topic vectors for each video event, where the weights reflected how reliably
270 each event was recalled. To visualize the result, we created a "wordle" image (Mueller et al., 2018)
271 where words weighted more heavily by better-remembered topics appear in a larger font (Fig. 6B,
272 green box). Events that reflected topics weighting heavily on characters like "Sherlock" and "John"
273 (i.e., the main characters) and locations like "221b Baker Street" (i.e., a major recurring location and
274 the address of the flat that Sherlock and John share) were best remembered. An analogous analysis
275 revealed which themes were poorly remembered. Here in computing the weighted average over
276 events' topic vectors we weighted each event in *inverse* proportion to how well it was remembered
277 (Fig. 6B, red box). This revealed that events with relatively minor characters such as "Mike,"
278 "Jeffrey," and "Molly," as well as less-integral plot locations (e.g., "hospital" and "office") were
279 least well-remembered. This suggests that what is retained in memory are the major plot elements
280 (i.e., the overall shape of what happened), whereas the more minor details are prone to pruning.

281 In addition to constructing overall summaries, assessing the video and recall topic vectors from
282 individual recalls can provide further insights. Specifically, for any given event we can construct
283 two wordles: one from the original video event's topic vector, and a second from the average topic

vectors produced by all participants' recalls of that event. We can then examine those wordles visually to gain an intuition for which aspects of the video event were recapitulated in participants' recalls of that event. Several example wordles are displayed in Figure 6C (wordles from the three best-remembered events are circled in green; wordles from the three worst-remembered events are circled in red). Using wordles to visually compare the topical content of each video event and the (average) corresponding recall event reveals the specific content from the specific events that is reliably retained in the transformation into memory (green events) or not (red events).

The results thus far inform us about which aspects of the dynamic content in the episode participants watched were preserved or altered in participants' memories of the episode. We next carried out a series of analyses aimed at understanding which brain structures might implement these processes. In one analysis we sought to identify which brain structures were sensitive to the video's dynamic content, as characterized by its topic trajectory. Specifically, we used a searchlight procedure to identify the extent to which each cluster of voxels exhibited a timecourse (as the participants watched the video) whose temporal correlation matrix matched the temporal correlation matrix of the original video's topic proportion matrix (Fig. 2B). As shown in Figure 7A, the analysis revealed a network of regions including bilateral frontal cortex and cingulate cortex, suggesting that these regions may play a role in maintaining information relevant to the narrative structure of the video. In a second analysis, we sought to identify which brain structures' responses (while viewing the video) reflected how each participant would later *recall* the video. We used an analogous searchlight procedure to identify clusters of voxels whose temporal correlation matrices reflected the temporal correlation matrix of the topic proportions for each individual's recalls (Figs. 2D, S4). As shown in Figure 7B, the analysis revealed a network of regions including the ventromedial prefrontal cortex (vmPFC), anterior cingulate cortex, and right medial temporal lobe (rMTL), suggesting that these regions may play a role in transforming each individual's experience into memory. In identifying regions whose responses to ongoing experiences reflect how those experiences will be remembered later, this latter analysis extends classic *subsequent memory analyses* (e.g., Paller and Wagner, 2002) to domain of naturalistic stimuli.

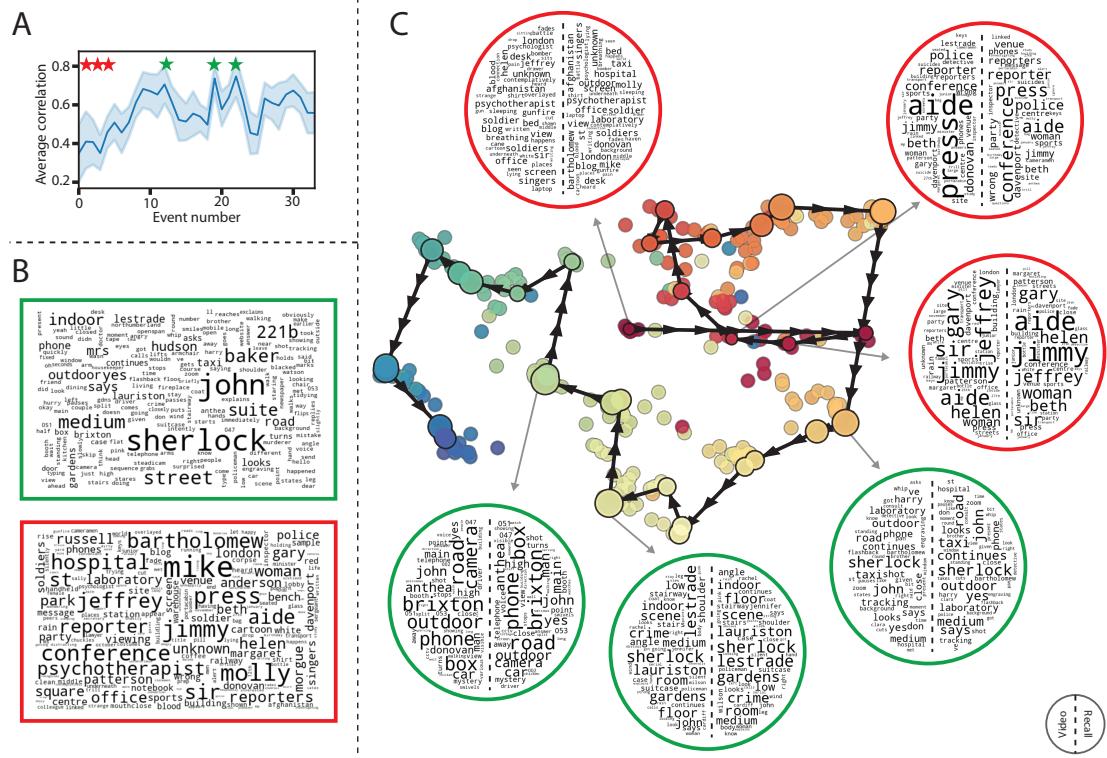


Figure 6: Transforming experience into memory. **A.** Average correlations (across participants) between the topic vectors from each video event and the closest-matching recall events. Error bars denote bootstrap-derived across-participant 95% confidence intervals. The stars denote the three best-remembered events (green) and worst-remembered events (red). **B.** Wordles comprising the top 200 highest-weighted words reflected in the weighted-average topic vector across video events. Green: video events were weighted by how well the topic vectors derived from recalls of those events matched the video events' topic vectors (Panel A). Red: video events were weighted by the inverse of how well their topic vectors matched the recalled topic vectors. **C.** The set of all video and recall events is projected onto the two-dimensional space derived in Figure 5. The dots outlined in black denote video events (dot size reflects the average correlation between the video event's topic vector and the topic vectors from the closest matching recalled events from each participant; bigger dots denote stronger correlations). The dots without black outlines denote recalled events. All dots are colored using the same scheme as Figure 5A. Wordles for several example events are displayed (green: three best-remembered events; red: three worst-remembered events). Within each circular wordle, the left side displays words associated with the topic vector for the video event, and the right side displays words associated with the (average) recall event topic vector, across all recall events matched to the given video event.

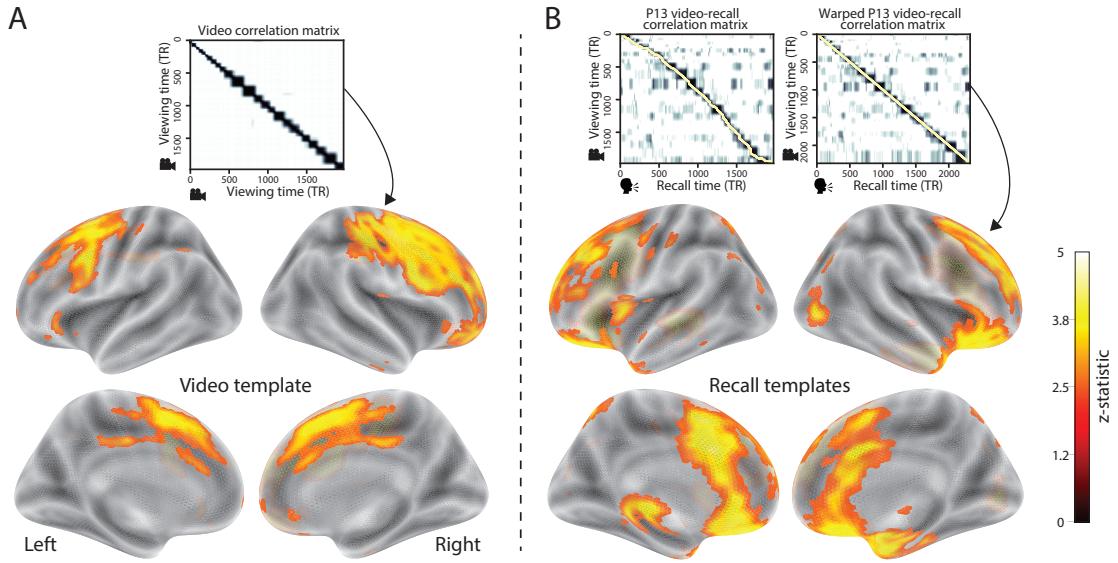


Figure 7: Brain structures that underlie the transformation of experience into memory. **A.** We searched for regions whose responses (as participants watched the video) matched the temporal correlation matrix of the video topic proportions. These regions are sensitive to the narrative structure of the video. **B.** We searched for regions whose responses (as participants watched the video) matched the temporal correlation matrix of the topic proportions derived from each individual's later recall of video. These regions are sensitive to how the narrative structure of the video is transformed into a memory of the video. Both panels: the maps are thresholded at $p < 0.05$, corrected.

311 **Discussion**

312 Our work casts remembering as reproducing (behaviorally and neurally) the topic trajectory, or
313 shape, of the original experience. This view draws inspiration from prior work aimed at elucidating
314 the neural and behavioral underpinnings of how we process dynamic naturalistic experiences
315 and remember them later. One approach to identifying neural responses to naturalistic stimuli
316 (including experiences) entails building a model of the stimulus and searching for brain regions
317 whose responses are consistent with the model. In prior work, a series of studies from Uri
318 Hasson's group (Lerner et al., 2011; Simony et al., 2016; Chen et al., 2017; Baldassano et al., 2017;
319 Zadbood et al., 2017) have extended this approach with a clever twist. Rather than building an
320 explicit stimulus model, these studies instead search for brain responses (while experiencing the
321 stimulus) that are reliably similar across individuals. So called *inter-subject correlation* (ISC) and
322 *inter-subject functional connectivity* (ISFC) analyses effectively treat other people's brain responses
323 to the stimulus as a "model" of how its features change over time. By contrast, in our present
324 work we used topic models and HMMs to construct an explicit stimulus model (i.e., the topic
325 trajectory of the video). When we searched for brain structures whose responses are consistent
326 with the video's topic trajectory, we identified a network of structures that overlapped strongly
327 with the "long temporal receptive window" network reported by the Hasson group (e.g., compare
328 our Fig. 7A with the map of long temporal receptive window voxels in Lerner et al., 2011). This
329 provides support for the notion that part of the long temporal receptive window network may be
330 maintaining an explicit model of the stimulus dynamics. When we performed a similar analysis
331 after swapping out the video's topic trajectory with the recall topic trajectories of each individual
332 participant, this allowed us to identify brain regions whose responses (as the participants viewed
333 the video) reflected how the video trajectory would be transformed in memory (as reflected by
334 the recall topic trajectories). The analysis revealed that the rMTL and vmPFC may play a role in
335 this person-specific transformation from experience into memory. The role of the MTL in episodic
336 memory encoding has been well-reported (e.g., Paller and Wagner, 2002; Davachi et al., 2003;
337 Ranganath et al., 2004; Davachi, 2006). Prior work has also implicated the medial prefrontal cortex

338 in representing “schema” knowledge (i.e., general knowledge about the format of an ongoing
339 experience given prior similar experiences; van Kesteren et al., 2012; Schlichting and Preston,
340 2015; Gilboa and Marlatte, 2017; Spalding et al., 2018). Integrating across our study and this prior
341 work, one interpretation is that the person-specific transformations mediated (or represented)
342 by the rMTL and vmPFC may reflect schema knowledge being leveraged, formed, or updated,
343 incorporating ongoing experience into previously acquired knowledge.

344 When modeling memory experiments, often times events (or items) and their later memories
345 are treated as binary (or categorical in the case of confidence ratings) and independent events.
346 Our novel framework allows one to assess memory performance in a more continuous way (e.g.
347 precision), as well as analyze the correlational structure of each encoding event to each memory
348 event (e.g. distinctiveness). Further and importantly, it allows for consideration of the actual
349 content of the experience/memories, which is not typically modeled. Leveraging this, using 2 novel
350 memory metrics we find that the successful memory performance is related to 1) the *precision* with
351 which the participant recounts each event and 2) how *distinctive* each recall event is (relative to the
352 other recalled events). The first finding suggests to us that the accuracy of recall for *any individual*
353 *event* may predict the overall amount of information recovered by the participant. The second
354 finding suggests that remembering/describing events in a unique way (relative to other recalled
355 events) is also related to the quantity of content recovered. Intriguingly, prior studies show that
356 pattern separation, or the ability to discriminate between similar experiences, is impaired in many
357 cognitive disorders as well as natural aging (Stark et al., 2010; Yassa et al., 2011; Yassa and Stark,
358 2011). Future work might explore how/whether the novel metrics introduce here compare between
359 cognitively impoverished groups and healthy controls.

360 While a large number of language models exist (e.g. WAS, LSA, word2vec, universal sentence
361 encoder) (Landauer et al., 1998; Cer et al., 2018), here we use topic models for a few reasons. First,
362 topic models capture the *essence* of a text passage devoid of the specific set and order of words
363 used. This was an important feature of our model since different people may accurately recall
364 a scene using very different language. Secondly, words can mean different things in different
365 contexts (e.g. baseball bat vs. the animal bat), and topic models are robust to this since words can

366 be part of multiple topics. Lastly, topic models provide a straight forward to recover the weights
367 for the particular words comprising a topic, allowing for easy interpretation of an event's contents
368 (e.g. Fig. 6). Other models such as Google's universal sentence encoder offer a context-sensitive
369 encoding of text passages, but the encoding space is complex and non-linear and thus, recovering
370 the original words used to fit the model is not straight forward. However, it's worth pointing out
371 that our framework is divorced from the particular choice of language model. Moreover, many of
372 the aspects of our framework could be swapped out for other choices. For example, the language
373 model, the timeseries segmentation model and the video-recall matching function could all be
374 customized for the particular problem. Indeed for some problems, recovery of the particular recall
375 words may not be necessary, and thus other text-modeling approaches (such as universal sentence
376 encoder) may be preferable. Future work will explore the influence of particular model choices on
377 the framework's accuracy.

378 Our work has broad implications for how we characterize and assess memory in real-world set-
379 tings such as the classroom or physician's office. For example, the most commonly used classroom
380 evaluation tools involve computing the proportion of correctly answered exam questions. Our
381 work indicates that this approach is only loosely related to what educators might really want to
382 measure: how well did the students understand the key ideas presented in the course? One could
383 apply the computational framework we developed to construct topic trajectories for the video and
384 participants' recalls to build explicit content models of the course material and exam questions.
385 This approach would provide a more nuanced and specific view into which aspects of the material
386 students had learned well (or poorly). In clinical settings, memory measures that incorporate such
387 explicit content models might also provide more direct evaluations of patients' memories.

388 **Methods**

389 **Experimental design and data collection**

390 Data were collected by Chen et al. (2017). In brief, participants ($n = 17$) viewed the first 48 minutes
391 of “A Study in Pink”, the first episode of the BBC television series *Sherlock*, while fMRI volumes
392 were collected (TR = 1500 ms). The stimulus was divided into a 23 min (946 TR) and a 25 min
393 (1030 TR) segment to mitigate technical issues related to the scanner. After finishing the clip,
394 participants were instructed to (quoting from Chen et al., 2017) “describe what they recalled of the
395 [episode] in as much detail as they could, to try to recount events in the original order they were
396 viewed in, and to speak for at least 10 minutes if possible but that longer was better. They were told
397 that completeness and detail were more important than temporal order, and that if at any point
398 they realized they had missed something, to return to it. Participants were then allowed to speak
399 for as long as they wished, and verbally indicated when they were finished (e.g., ‘I’m done’).”
400 For additional details about the experimental procedure and scanning parameters see Chen et al.
401 (2017). The experimental protocol was approved by Princeton University’s Institutional Review
402 Board.

403 After preprocessing the fMRI data and warping the images into a standard (3 mm³ MNI) space,
404 the voxel activations were z-scored (within voxel) and spatially smoothed using a 6 mm (full width
405 at half maximum) Gaussian kernel. The fMRI data were also cropped so that all video-viewing
406 data were aligned across participants. This included a constant 3 TR (4.5 s) shift to account for the
407 lag in the hemodynamic response. (All of these preprocessing steps followed Chen et al., 2017,
408 where additional details may be found.)

409 **Data and code availability**

410 The fMRI data we analyzed are available online [here](#). The behavioral data and all of our analysis
411 code may be downloaded [here](#).

412 **Statistics**

413 All statistical tests we performed were two-sided.

414 **Modeling the dynamic content of the video and recall transcripts**

415 **Topic modeling**

416 The input to the topic model we trained to characterize the dynamic content of the video comprised
417 hand-generated annotations of each of 1000 scenes spanning the video clip (generated by Chen
418 et al., 2017). The features included: narrative details (a sentence or two describing what happened
419 in that scene); whether the scene took place indoors or outdoors; names of any characters that
420 appeared in the scene; name(s) of characters in camera focus; name(s) of characters who were
421 speaking in the scene; the location (in the story) that the scene took place; camera angle (close
422 up, medium, long, top, tracking, over the shoulder, etc.); whether music was playing in the
423 scene or not; and a transcription of any on-screen text. We concatenated the text for all of these
424 features within each segment, creating a “bag of words” describing each scene. We then re-
425 organized the text descriptions into overlapping sliding windows spanning 50 scenes each. In
426 other words, the first text sample comprised the combined text from the first 50 scenes (i.e., 1–50),
427 the second comprised the text from scenes 2–51, and so on. We trained our model using these
428 overlapping text samples with `scikit-learn` (version 0.19.1; Pedregosa et al., 2011), called from
429 our high-dimensional visualization and text analysis software, `HyperTools` (Heusser et al., 2018b).
430 Specifically, we use the `CountVectorizer` class to transform the text from each scene into a vector of
431 word counts (using the union of all words across all scenes as the “vocabulary,” excluding English
432 stop words); this yields a number-of-scenes by number-of-words *word count* matrix. We then
433 use the `LatentDirichletAllocation` class (`topics=100, method='batch'`) to fit a topic model (Blei
434 et al., 2003) to the word count matrix, yielding a number-of-scenes (1000) by number-of-topics
435 (100) *topic proportions* matrix. The topic proportions matrix describes which mix of topics (latent
436 themes) is present in each scene. Next, we transformed the topic proportions matrix to match the
437 1976 fMRI volume acquisition times. For each fMRI volume, we took the topic proportions from

438 whatever scene was displayed for most of that volume's 1500 ms acquisition time. This yielded a
439 new number-of-TRs (1976) by number-of-topics (100) topic proportions matrix.

440 We created similar topic proportions matrices using hand-annotated transcripts of each participant
441's recall of the video (annotated by Chen et al., 2017). We tokenized the transcript into a list of
442 sentences, and then re-organized the list into overlapping sliding windows spanning 10 sentences
443 each; in turn we transformed each window's sentences into a word count vector (using the same
444 vocabulary as for the video model). We then used the topic model already trained on the video
445 scenes to compute the most probable topic proportions for each sliding window. This yielded a
446 number-of-sentences (range: 68–294) by number-of-topics (100) topic proportions matrix, for each
447 participant. These reflected the dynamic content of each participant's recalls. Note: for details
448 on how we selected the video and recall window lengths and number of topics, see *Supporting*
449 *Information* and Figure S1.

450 **Parsing topic trajectories into events using Hidden Markov Models**

451 We parsed the topic trajectories of the video and participants' recalls into events using Hidden
452 Markov Models (Rabiner, 1989). Given the topic proportions matrix (describing the mix of topics
453 at each timepoint) and a number of states, K , an HMM recovers the set of state transitions that
454 segments the timeseries into K discrete states. Following Baldassano et al. (2017), we imposed an
455 additional set of constraints on the discovered state transitions that ensured that each state was
456 encountered exactly once (i.e., never repeated). We used the BrainIAK toolbox (Capota et al., 2017)
457 to implement this segmentation.

458 We used an optimization procedure to select the appropriate K for each topic proportions
459 matrix. Specifically, we computed (for each matrix)

$$\operatorname{argmax}_K \left[\frac{a}{b} - \frac{K}{\alpha} \right],$$

460 where a was the average correlation between the topic vectors of timepoints within the same state;
461 b was the average correlation between the topic vectors of timepoints within *different* states; and

462 α was a regularization parameter that we set to 5 times the window length (i.e., 250 scenes for
463 the video topic trajectory and 50 sentences for the recall topic trajectories). Figure 2B displays the
464 event boundaries returned for the video, and Figure S4 displays the event boundaries returned
465 for each participant's recalls. After obtaining these event boundaries, we created stable estimates
466 of each topic proportions matrix by averaging the topic vectors within each event. This yielded a
467 number-of-events by number-of-topics matrix for the video and recalls from each participant.

468 We also evaluated a parameter-free procedure for choosing K , which finds the K value that
469 maximizes the Wasserstein distance (a.k.a. "Earth mover's" distance) between the within and
470 across event distributions of correlation values. This alternative procedure largely replicated the
471 pattern of results found with the parameterized method described above, but recovered sub-
472 stantially fewer events on average (Fig.S6). While both approaches seem to underestimate the
473 number of video/recall events relative to the "true" number (as determined by human raters), the
474 parameterized approach was closer to the true number.

475 **Naturalistic extensions of classic list-learning analyses**

476 In traditional list-learning experiments, participants view a list of items (e.g., words) and then recall
477 the items later. Our video-recall event matching approach affords us the ability to analyze memory
478 in a similar way. The video and recall events can be treated analogously to studied and recalled
479 "items" in a list-learning study. We can then extend classic analyses of memory performance and
480 dynamics (originally designed for list-learning experiments) to the more naturalistic video recall
481 task used in our study.

482 Perhaps the simplest and most widely used measure of memory performance is *accuracy*— i.e.,
483 the proportion of studied (experienced) items (in this case, the 34 video events) that the participant
484 later remembered. Chen et al. (2017) developed a human rating system whereby the quality of
485 each participant's memory was evaluated by an independent rater. We found a strong across-
486 participants correlation between these independant ratings and the overall number of events that
487 our HMM approach identified in participants' recalls (Pearson's $r(15) = 0.67, p = 0.003$).

488 As described below, we next considered a number of memory performance measures that are

489 typically associated with list-learning studies. We also provide a software package, Quail, for
490 carrying out these analyses (Heusser et al., 2017).

491 **Probability of first recall (PFR).** PFR curves (Welch and Burnett, 1924; Postman and Phillips,
492 1965; Atkinson and Shiffrin, 1968) reflect the probability that an item will be recalled first as a
493 function of its serial position during encoding. To carry out this analysis, we initialized a number-
494 of-participants (17) by number-of-video-events (34) matrix of zeros. Then for each participant, we
495 found the index of the video event that was recalled first (i.e., the video event whose topic vector
496 was most strongly correlated with that of the first recall event) and filled in that index in the matrix
497 with a 1. Finally, we averaged over the rows of the matrix, resulting in a 1 by 34 array representing
498 the proportion of participants that recalled an event first, as a function of the order of the event's
499 appearance in the video (Fig. 3A).

500 **Lag conditional probability curve (lag-CRP).** The lag-CRP curve (Kahana, 1996) reflects the
501 probability of recalling a given event after the just-recalled event, as a function of their relative
502 positions (or *lag*). In other words, a lag of 1 indicates that a recalled event came immediately after
503 the previously recalled event in the video, and a lag of -3 indicates that a recalled event came 3
504 events before the previously recalled event. For each recall transition (following the first recall),
505 we computed the lag between the current recall event and the next recall event, normalizing by
506 the total number of possible transitions. This yielded a number-of-participants (17) by number-
507 of-lags (-33 to +33; 67 lags total) matrix. We averaged over the rows of this matrix to obtain a
508 group-averaged lag-CRP curve (Fig. 3B).

509 **Serial position curve (SPC).** SPCs (Murdock, 1962) reflect the proportion of participants that
510 remember each item as a function of their serial position during encoding. We initialized a number-
511 of-participants (17) by number-of-video-events (34) matrix of zeros. Then, for each recalled event,
512 for each participant, we found the index of the video event that the recalled event most closely
513 matched (via the correlation between the events' topic vectors) and entered a 1 into that position
514 in the matrix (i.e., for the given participant and event). This resulted in a matrix whose entries

515 indicated whether or not each event was recalled by each participant (depending on whether the
516 corresponding entires were set to one or zero). Finally, we averaged over the rows of the matrix
517 to yield a 1 by 34 array representing the proportion of participants that recalled each event as a
518 function of the order of the event's appearance in the video (Fig. 3C).

519 **Temporal clustering scores.** Temporal clustering refers to the extent to which participants group
520 their recall responses according to encoding position (Polyn et al., 2009). For instance, if a par-
521 ticipant recalled the video events in the exact order they occurred (or in exact reverse order), this
522 would yield a score of 1. If a participant recalled the events in random order, this would yield
523 an expected score of 0.5. For each recall event transition (and separately for each participant), we
524 sorted all not-yet-recalled events according to their absolute lag (i.e., distance away in the video).
525 We then computed the percentile rank of the next event the participant recalled. We averaged
526 these percentile ranks across all of the participant's recalls to obtain a single temporal clustering
527 score for the participant (mean: 0.808, SEM: 0.022). Overall, we found that participants with higher
528 temporal clustering scores also tended to recall more events (Pearson's $r(15) = 0.62, p = 0.007$).

529 **Semantic clustering scores.** Semantic clustering measures the extent to which participants clus-
530 tered their recall responses according to semantic similarity (Polyn et al., 2009). Here, we used the
531 topic vectors for each event as a proxy for its semantic content. Thus, the similarity between the
532 semantic content for two events can be computed by correlating their respective topic vectors. For
533 each recall event transition, we sorted all not-yet-recalled events according to how correlated the
534 topic vector of *the closest-matching video event* was to the topic vector of the closest-matching video
535 event to the just-recalled event. We then computed the percentile rank of the observed next recall.
536 We averaged these percentile ranks across all of the participant's recalls to obtain a single semantic
537 clustering score for the participant (mean: 0.813, SEM: 0.022). We found that participants who
538 exhibited stronger semantic clustering scores overall remembered more video events (Pearson's
539 $r(15) = 0.55, p = 0.02$).

540 **Novel naturalistic memory metrics**

541 **Precision.** We tested whether participants who recalled more events were also more *precise* in their
542 recollections. For each participant, we computed the correlation between the topic vectors for each
543 recall event and that of its closest-matching video event (only for the events which they recalled).
544 We Fisher's z-transformed the correlations, computed the average and then inverse Fisher's z-
545 transformed the resulting value. This gave a single value per participant representing the average
546 precision across all recalled events. We then correlated this value with hand-annotated as well as
547 model derived (e.g. k or the number of events recovered by the HMM) memory performance.

548 **Distinctiveness.** We also considered the *distinctiveness* of each recalled event. That is, how
549 uniquely a recalled event's topic vector matched a given video event topic vector, versus the
550 topic vectors for the other video events. We hypothesized that participants with high memory
551 performance might describe each event in a more distinctive way (relative to those with lower
552 memory performance who might describe events in a more general way). To test this hypothesis
553 we define a distinctiveness score for each recalled event as

$$d(\text{event}) = 1 - \bar{c}(\text{event}),$$

554 where $\bar{c}(\text{event})$ is the average correlation between the given recalled event's topic vector and the
555 topic vectors from all video events *except* the best-matching video event. We then averaged these
556 distinctiveness scores across all of the events recalled by the given participant. As above, we used
557 Fisher's z (transform and inverse-transform) before/after averaging correlation values. Finally,
558 we correlated these values with hand-annotated and model derived memory performance scores
559 across-subjects.

560 **Visualizing the video and recall topic trajectories**

561 We used the UMAP algorithm (McInnes and Healy, 2018) to project the 100-dimensional topic space
562 onto a two-dimensional space for visualization (Figs. 5, 6). To ensure that all of the trajectories were

563 projected onto the *same* lower dimensional space, we computed the low-dimensional embedding
564 on a “stacked” matrix created by vertically concatenating the events-by-topics topic proportions
565 matrices for the video and all 17 participants’ recalls. We then divided the rows of the result (a
566 total-number-of-events by two matrix) back into separate matrices for the video topic trajectory
567 and the trajectories for each participant’s recalls (Fig. 5). This general approach for discovering
568 a shared low-dimensional embedding for a collections of high-dimensional observations follows
569 Heusser et al. (2018b).

570 **Estimating the consistency of flow through topic space across participants**

571 In Figure 5B, we present an analysis aimed at characterizing locations in topic space that dif-
572 ferent participants move through in a consistent way (via their recall topic trajectories). The
573 two-dimensional topic space used in our visualizations (Fig. 5) ranged from -5 to 5 (arbitrary) units
574 in the x dimension and from -6.5 to 2 units in the y dimension. We divided this space into a grid
575 of vertices spaced 0.25 units apart. For each vertex, we examined the set of line segments formed
576 by connecting each pair successively recalled events, across all participants, that passed within 0.5
577 units. We computed the distribution of angles formed by those segments and the x -axis, and used a
578 Rayleigh test to determine whether the distribution of angles was reliably “peaked” (i.e., consistent
579 across all transitions that passed through that local portion of topic space). To create Figure 5B we
580 drew an arrow originating from each grid vertex, pointing in the direction of the average angle
581 formed by line segments that passed within 0.5 units. We set the arrow lengths to be inversely
582 proportional to the p -values of the Rayleigh tests at each vertex. Specifically, for each vertex we
583 converted all of the angles of segments that passed within 0.5 units to unit vectors, and we set
584 the arrow lengths at each vertex proportional to the length of the (circular) mean vector. We also
585 indicated any significant results ($p < 0.05$, corrected using the Benjamani-Hochberg procedure) by
586 coloring the arrows in blue (darker blue denotes a lower p -value, i.e., a longer mean vector); all
587 tests with $p \geq 0.05$ are displayed in gray and given a lower opacity value.

588 **Searchlight fMRI analyses**

589 In Figure 7, we present two analyses aimed at identifying brain structures whose responses (as
590 participants viewed the video) exhibited particular temporal correlations. We developed a search-
591 light analysis whereby we constructed a cube centered on each voxel (radius: 5 voxels). For each
592 of these cubes, we computed the temporal correlation matrix of the voxel responses during video
593 viewing. Specifically, for each of the 1976 volumes collected during video viewing, we correlated
594 the activity patterns in the given cube with the activity patterns (in the same cube) collected during
595 every other timepoint. This yielded a 1976 by 1976 correlation matrix for each cube.

596 Next, we constructed two sets of “template” matrices: one reflected the video’s topic trajectory
597 and the other reflected each participant’s recall topic trajectory. To construct the video template, we
598 computed the correlations between the topic proportions estimated for every pair of TRs (prior to
599 segmenting the trajectory into discrete events; i.e., the correlation matrix shown in Figs. 2B and 7A).
600 We constructed similar temporal correlation matrices for each participant’s recall topic trajectory
601 (Figs. 2D, S4). However, to correct for length differences and potential non-linear transformations
602 between viewing time and recall time, we first used dynamic time warping (Berndt and Clifford,
603 1994) to temporally align participants’ recall topic trajectories with the video topic trajectory (an
604 example correlation matrix before and after warping is shown in Fig. 7B). This yielded a 1976 by
605 1976 correlation matrix for the video template and for each participant’s recall template.

606 To determine which (cubes of) voxel responses reliably matched the video template, we cor-
607 related the upper triangle of the voxel correlation matrix for each cube with the upper triangle
608 of the video template matrix (Kriegeskorte et al., 2008). This yielded, for each participant, a
609 single correlation value. We computed the average (Fisher z-transformed) correlation coefficient
610 across participants. We used a permutation-based procedure to assess significance, whereby we
611 re-computed the average correlations for each of 100 “null” video templates (constructed by circu-
612 larly shifting the template by a random number of timepoints). (For each permutation, the same
613 shift was used for all participants.) We then estimated a p -value by computing the proportion of
614 shifted correlations that were larger than the observed (unshifted) correlation. To create the map

615 in Figure 7A we thresholded out any voxels whose correlation values fell below the 95th percentile
616 of the permutation-derived null distribution.

617 We used a similar procedure to identify which voxels' responses reflected the recall templates.
618 For each participant, we correlated the upper triangle of the correlation matrix for each cube of
619 voxels with their (time warped) recall correlation matrix. As in the video template analysis this
620 yielded a single correlation coefficient for each participant. However, whereas the video analysis
621 compared every participant's responses to the same template, here the recall templates were
622 unique for each participant. We computed the average z -transformed correlation coefficient across
623 participants, and used the same permutation procedure we developed for the video responses to
624 assess significant correlations. To create the map in Figure 7B we thresholded out any voxels whose
625 correlation values fell below the 95th percentile of the permutation-derived null distribution.

626 References

- 627 Atkinson, R. C. and Shiffrin, R. M. (1968). Human memory: A proposed system and its control
628 processes. In Spence, K. W. and Spence, J. T., editors, *The psychology of learning and motivation*,
629 volume 2, pages 89–105. Academic Press, New York.
- 630 Baldassano, C., Chen, J., Zadbood, A., Pillow, J. W., Hasson, U., and Norman, K. A. (2017).
631 Discovering event structure in continuous narrative perception and memory. *Neuron*, 95(3):709–
632 721.
- 633 Berndt, D. J. and Clifford, J. (1994). Using dynamic time warping to find patterns in time series. In
634 *KDD workshop*, volume 10, pages 359–370.
- 635 Blei, D. M. and Lafferty, J. D. (2006). Dynamic topic models. In *Proceedings of the 23rd International
636 Conference on Machine Learning*, ICML '06, pages 113–120, New York, NY, US. ACM.
- 637 Blei, D. M., Ng, A. Y., and Jordan, M. I. (2003). Latent dirichlet allocation. *Journal of Machine
638 Learning Research*, 3:993 – 1022.

- 639 Brunec, I. K., Moscovitch, M. M., and Barense, M. D. (2018). Boundaries shape cognitive representations of spaces and events. *Trends in Cognitive Sciences*, 22(7):637–650.
- 640
- 641 Capota, M., Turek, J., Chen, P.-H., Zhu, X., Manning, J. R., Sundaram, N., Keller, B., Wang, Y., and
642 Shin, Y. S. (2017). Brain imaging analysis kit.
- 643 Cer, D., Yang, Y., y Kong, S., Hua, N., Limtiaco, N., John, R. S., Constant, N., Guajardo-Cespedes,
644 M., Yuan, S., Tar, C., Sung, Y.-H., Strope, B., and Kurzweil, R. (2018). Universal sentence encoder.
645 *arXiv*, 1803.11175.
- 646 Chen, J., Leong, Y. C., Honey, C. J., Yong, C. H., Norman, K. A., and Hasson, U. (2017). Shared
647 memories reveal shared structure in neural activity across individuals. *Nature Neuroscience*,
648 20(1):115.
- 649 Clewett, D. and Davachi, L. (2017). The ebb and flow of experience determines the temporal
650 structure of memory. *Curr Opin Behav Sci*, 17:186–193.
- 651 Davachi, L. (2006). Item, context and relational episodic encoding in humans. *Current Opinion in*
652 *Neurobiology*, 16(6):693—700.
- 653 Davachi, L., Mitchell, J. P., and Wagner, A. D. (2003). Multiple routes to memory: distinct medial
654 temporal lobe processes build item and source memories. *Proceedings of the National Academy of*
655 *Sciences, USA*, 100(4):2157 – 2162.
- 656 DuBrow, S. and Davachi, L. (2013). The influence of contextual boundaries on memory for the
657 sequential order of events. *Journal of Experimental Psychology: General*, 142(4):1277–1286.
- 658 Ezzyat, Y. and Davachi, L. (2011). What constitutes an episode in episodic memory? *Psychological*
659 *Science*, 22(2):243–252.
- 660 Gilboa, A. and Marlatte, H. (2017). Neurobiology of schemas and schema-mediated memory.
661 *Trends Cogn Sci*, 21(8):618–631.

- 662 Heusser, A. C., Ezzyat, Y., Shiff, I., and Davachi, L. (2018a). Perceptual boundaries cause mnemonic
663 trade-offs between local boundary processing and across-trial associative binding. *Journal of*
664 *Experimental Psychology Learning, Memory, and Cognition*, 44(7):1075–1090.
- 665 Heusser, A. C., Fitzpatrick, P. C., Field, C. E., Ziman, K., and Manning, J. R. (2017). Quail: a
666 Python toolbox for analyzing and plotting free recall data. *The Journal of Open Source Software*,
667 10.21105/joss.00424.
- 668 Heusser, A. C., Ziman, K., Owen, L. L. W., and Manning, J. R. (2018b). HyperTools: a Python
669 toolbox for gaining geometric insights into high-dimensional data. *Journal of Machine Learning*
670 *Research*, 18(152):1–6.
- 671 Howard, M. W. and Kahana, M. J. (2002). A distributed representation of temporal context. *Journal*
672 *of Mathematical Psychology*, 46:269–299.
- 673 Howard, M. W., MacDonald, C. J., Tiganj, Z., Shankar, K. H., Du, Q., Hasselmo, M. E., and H., E.
674 (2014). A unified mathematical framework for coding time, space, and sequences in the medial
675 temporal lobe. *Journal of Neuroscience*, 34(13):4692–4707.
- 676 Howard, M. W., Viskontas, I. V., Shankar, K. H., and Fried, I. (2012). Ensembles of human MTL
677 neurons “jump back in time” in response to a repeated stimulus. *Hippocampus*, 22:1833–1847.
- 678 Huk, A., Bonnen, K., and He, B. J. (2018). Beyond trial-based paradigms: continuous behavior, on-
679 going neural activity, and naturalistic stimuli. *Journal of Neuroscience*, 10.1523/JNEUROSCI.1920-
680 17.2018.
- 681 Kahana, M. J. (1996). Associative retrieval processes in free recall. *Memory & Cognition*, 24:103–109.
- 682 Koriat, A. and Goldsmith, M. (1994). Memory in naturalistic and laboratory contexts: distin-
683 guishing accuracy-oriented and quantity-oriented approaches to memory assessment. *Journal of*
684 *Experimental Psychology: General*, 123(3):297–315.
- 685 Kriegeskorte, N., Mur, M., and Bandettini, P. (2008). Representational similarity analysis – con-
686 necting the branches of systems neuroscience. *Frontiers in Systems Neuroscience*, 2:1 – 28.

- 687 Landauer, T. K., Foltz, P. W., and Laham, D. (1998). Introduction to latent semantic analysis.
- 688 *Discourse Processes*, 25:259–284.
- 689 Lerner, Y., Honey, C. J., Silbert, L. J., and Hasson, U. (2011). Topographic mapping of a hierarchy
- 690 of temporal receptive windows using a narrated story. *Journal of Neuroscience*, 31(8):2906–2915.
- 691 Manning, J. (2019). Episodic memory: mental time travel or a quantum ‘memory wave’ function?
- 692 *PsyArXiv*, doi:10.31234/osf.io/6zjwb.
- 693 Manning, J. R., Norman, K. A., and Kahana, M. J. (2015). The role of context in episodic memory.
- 694 In Gazzaniga, M., editor, *The Cognitive Neurosciences, Fifth edition*, pages 557–566. MIT Press.
- 695 Manning, J. R., Polyn, S. M., Baltuch, G., Litt, B., and Kahana, M. J. (2011). Oscillatory patterns
- 696 in temporal lobe reveal context reinstatement during memory search. *Proceedings of the National*
- 697 *Academy of Sciences, USA*, 108(31):12893–12897.
- 698 McInnes, L. and Healy, J. (2018). UMAP: Uniform manifold approximation and projection for
- 699 dimension reduction. *arXiv*, 1802(03426).
- 700 Mueller, A., Fillion-Robin, J.-C., Boidol, R., Tian, F., Nechifor, P., yoonsubKim, Peter, Rampin, R.,
- 701 Corvellec, M., Medina, J., Dai, Y., Petrushev, B., Langner, K. M., Hong, Alessio, Ozsvald, I.,
- 702 vkolmakov, Jones, T., Bailey, E., Rho, V., IgorAPM, Roy, D., May, C., foobuzz, Piyush, Seong,
- 703 L. K., Goey, J. V., Smith, J. S., Gus, and Mai, F. (2018). WordCloud 1.5.0: a little word cloud
- 704 generator in Python. *Zenodo*, <https://zenodo.org/record/1322068#.W4tPKZNKh24>.
- 705 Murdock, B. B. (1962). The serial position effect of free recall. *Journal of Experimental Psychology*,
- 706 64:482–488.
- 707 Paller, K. A. and Wagner, A. D. (2002). Observing the transformation of experience into memory.
- 708 *Trends in Cognitive Sciences*, 6(2):93–102.
- 709 Pedregosa, F., Varoquaux, G., Gramfort, A., Michel, V., Thirion, B., Grisel, O., Blondel, M., Pretten-
- 710 hofer, P., Weiss, R., Dubourg, V., Vanderplas, J., Passos, A., Cournapeau, D., Brucher, M., Perrot,

- 711 M., and Duchesnay, E. (2011). Scikit-learn: Machine learning in Python. *Journal of Machine*
712 *Learning Research*, 12:2825–2830.
- 713 Polyn, S. M., Norman, K. A., and Kahana, M. J. (2009). A context maintenance and retrieval model
714 of organizational processes in free recall. *Psychological Review*, 116(1):129–156.
- 715 Postman, L. and Phillips, L. W. (1965). Short-term temporal changes in free recall. *Quarterly Journal*
716 *of Experimental Psychology*, 17:132–138.
- 717 Rabiner, L. (1989). A tutorial on Hidden Markov Models and selected applications in speech
718 recognition. *Proceedings of the IEEE*, 77(2):257–286.
- 719 Radvansky, G. A. and Zacks, J. M. (2017). Event boundaries in memory and cognition. *Curr Opin*
720 *Behav Sci*, 17:133–140.
- 721 Ranganath, C., Cohen, M. X., Dam, C., and D’Esposito, M. (2004). Inferior temporal, prefrontal,
722 and hippocampal contributions to visual working memory maintenance and associative memory
723 retrieval. *Journal of Neuroscience*, 24(16):3917–3925.
- 724 Ranganath, C. and Ritchey, M. (2012). Two cortical systems for memory-guided behavior. *Nature*
725 *Reviews Neuroscience*, 13:713 – 726.
- 726 Schlichting, M. L. and Preston, A. R. (2015). Memory integration: neural mechanisms and impli-
727 cations for behavior. *Current Opinion in Behavioral Sciences*, 1:1–8.
- 728 Simony, E., Honey, C. J., Chen, J., and Hasson, U. (2016). Uncovering stimulus-locked network
729 dynamics during narrative comprehension. *Nature Communications*, 7(12141):1–13.
- 730 Spalding, K. N., Schlichting, M. L., Zeithamova, D., Preston, A. R., Tranel, D., Duff, M. C., and
731 Warren, D. E. (2018). Ventromedial prefrontal cortex is necessary for normal associative inference
732 and memory integration. *The Journal of Neuroscience*, 38(15):3767–3775.
- 733 Stark, S. M., Yassa, M. A., and Stark, C. E. L. (2010). Individual differences in spatial pattern
734 separation performance associated with healthy aging in humans. *Learning & Memory*, 17(6):284–
735 288.

- 736 Tomrary, A. and Davachi, L. (2017). Consolidation promotes the emergence of representational
737 overlap in the hippocampus and medial prefrontal cortex. *Neuron*, 96(1):228–241.
- 738 van Kesteren, M. T. R., Ruiter, D. J., Fernández, G., and Henson, R. N. (2012). How schema and
739 novelty augment memory formation. *Trends Neurosci*, 35(4):211–9.
- 740 Waskom, M., Botvinnik, O., Okane, D., Hobson, P., David, Halchenko, Y., Lukauskas, S., Cole, J. B.,
741 Warmenhoven, J., de Ruiter, J., Hoyer, S., Vanderplas, J., Villalba, S., Kunter, G., Quintero, E.,
742 Martin, M., Miles, A., Meyer, K., Augspurger, T., Yarkoni, T., Bachant, P., Williams, M., Evans,
743 C., Fitzgerald, C., Brian, Wehner, D., Hitz, G., Ziegler, E., Qalieh, A., and Lee, A. (2016). Seaborn:
744 v0.7.1.
- 745 Welch, G. B. and Burnett, C. T. (1924). Is primacy a factor in association-formation. *American Journal
746 of Psychology*, 35:396–401.
- 747 Yassa, M. A., Lacy, J. W., Stark, S. M., Albert, M. S., Gallagher, M., and Stark, C. E. L. (2011). Pattern
748 separation deficits associated with increased hippocampal ca3 and dentate gyrus activity in
749 nondemented older adults. *Hippocampus*, 21(9):968–979.
- 750 Yassa, M. A. and Stark, C. E. L. (2011). Pattern separation in the hippocampus. *Trends In Neuro-
751 sciences*, 34(10):515–525.
- 752 Yonelinas, A. P. (2002). The nature of recollection and familiarity: A review of 30 years of research.
753 *Journal of Memory and Language*, 46:441–517.
- 754 Zacks, J. M., Speer, N. K., Swallow, K. M., Braver, T. S., and Reynolds, J. R. (2007). Event perception:
755 a mind-brain perspective. *Psychological Bulletin*, 133:273–293.
- 756 Zadbood, A., Chen, J., Leong, Y. C., Norman, K. A., and Hasson, U. (2017). How we transmit
757 memories to other brains: Constructing shared neural representations via communication. *Cereb
758 Cortex*, 27(10):4988–5000.
- 759 Zwaan, R. A. and Radvansky, G. A. (1998). Situation models in language comprehension and
760 memory. *Psychological Bulletin*, 123(2):162 – 185.

⁷⁶¹ **Supporting information**

⁷⁶² Supporting information is available in the online version of the paper.

⁷⁶³ **Acknowledgements**

⁷⁶⁴ We thank Luke Chang, Janice Chen, Chris Honey, Lucy Owen, Emily Whitaker, and Kirsten Ziman
⁷⁶⁵ for feedback and scientific discussions. We also thank Janice Chen, Yuan Chang Leong, Kenneth
⁷⁶⁶ Norman, and Uri Hasson for sharing the data used in our study. Our work was supported in part
⁷⁶⁷ by NSF EPSCoR Award Number 1632738. The content is solely the responsibility of the authors
⁷⁶⁸ and does not necessarily represent the official views of our supporting organizations.

⁷⁶⁹ **Author contributions**

⁷⁷⁰ Conceptualization: A.C.H. and J.R.M.; Methodology: A.C.H. and J.R.M.; Software: A.C.H., P.C.F.
⁷⁷¹ and J.R.M.; Analysis: A.C.H., P.C.F. and J.R.M.; Writing, Reviewing, and Editing: A.C.H., P.C.F.
⁷⁷² and J.R.M.; Supervision: J.R.M.

⁷⁷³ **Author information**

⁷⁷⁴ The authors declare no competing financial interests. Correspondence and requests for materials
⁷⁷⁵ should be addressed to J.R.M. (jeremy.r.manning@dartmouth.edu).