

¹ Memory for television episodes preserves event content
² while introducing new across-event similarities

³ Andrew C. Heusser^{1,2}, Paxton C. Fitzpatrick¹, and Jeremy R. Manning^{1,*}

¹Department of Psychological and Brain Sciences

Dartmouth College, Hanover, NH 03755, USA

²Akili Interactive

Boston, MA 02110

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

⁴ December 11, 2019

⁵ **Abstract**

The ways our experiences unfold over time define unique *trajectories* through the relevant representational spaces. Within this geometric framework, one can compare the shape of the trajectory formed by an experience to that defined by our later remembering of that experience. We propose a framework for mapping naturalistic experiences onto geometric spaces that characterize how experiences are segmented into discrete events, and how the contents of event sequences evolve over time. We apply this approach to a naturalistic memory experiment which had participants view and recount a television episode. The content of participants' recounts of events from the original episode closely matched the original episode's content. However, the similarity patterns *across* events was much different in the original episode as compared with participants' recounts. We also identified a network of brain structures that are sensitive to the "shapes" of ongoing experiences, and an overlapping network that is sensitive (at the time of encoding) to how people later remembered those experiences in relation to other experiences.

18 In this way, modeling the content of richly structured experiences can reveal how (geometrically
19 and conceptually) those experiences are segmented into events and integrated into our memories
20 of other experiences.

21 **Introduction**

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

42 How might one go about formally characterizing the “essence” of an experience, or whether
43 it has been recovered by the rememberer? Any given moment of an experience derives meaning

44 from surrounding moments, as well as from longer-range temporal associations (??). Therefore,
45 the timecourse describing how an event unfolds is fundamental to its overall meaning. Further,
46 this hierarchy formed by our subjective experiences at different timescales defines a *context* for
47 each new moment (e.g., ??), and plays an important role in how we interpret that moment and
48 remember it later (for review see ?). Our memory systems can leverage these associations to form
49 predictions that help guide our behaviors (?). For example, as we navigate the world, the features
50 of our subjective experiences tend to change gradually (e.g., the room or situation we are in at any
51 given moment is strongly temporally autocorrelated), allowing us to form stable estimates of our
52 current situation and behave accordingly (??).

53 Although our experiences most often change gradually, they also occasionally change sud-
54 denly (e.g., when we walk through a doorway; ?). Prior research suggests that these sharp
55 transitions (termed *event boundaries*) during an experience help to discretize our experiences (and
56 their mental representations) into *events* (??????). The interplay between the stable (within event)
57 and transient (across event) temporal dynamics of an experience also provides a potential frame-
58 work for transforming experiences into memories that distill those experiences down to their
59 essence. For example, prior work has shown that event boundaries can influence how we learn
60 sequences of items (??), navigate (?), and remember and understand narratives (??). Prior research
61 has implicated the hippocampus and the medial prefrontal cortex as playing a critical role in
62 transforming experiences into stuctured and consolidated memories (?).

63 Here we sought to examine how the temporal dynamics of a “naturalistic” experience were
64 later reflected in participants’ memories. We analyzed an open dataset that comprised behavioral
65 and functional Magnetic Resonance Imaging (fMRI) data collected as participants viewed and
66 then verbally recounted an episode of the BBC television series *Sherlock* (?). We developed a
67 computational framework for characterizing the temporal dynamics of the moment-by-moment
68 content of the episode and of participants’ verbal recalls. Specifically, we use topic modeling (?) to
69 characterize the thematic conceptual (semantic) content present in each moment of the episode and
70 recalls, and Hidden Markov Models (??) to discretize this evolving semantic content into events.
71 In this way, we cast naturalistic experiences (and recalls of those experiences) as *trajectories* that

72 describe how the experiences evolve over time. Under this framework, successful remembering
73 entails verbally “traversing” the content trajectory of the episode, thereby reproducing the shape
74 (or essence) of the original experience. Comparing the shapes of the topic trajectories of the
75 episode and of participants’ retellings of the episode then reveals which aspects of the episode
76 were preserved (or lost) in the translation into memory. We further examine whether 1) the
77 *precision* with which a participant recounts each event and 2) the *distinctiveness* each recall event is
78 (relative to the other recalled events) relates to their overall memory performance. Last, we identify
79 networks of brain structures whose responses (as participants watched the episode) reflected the
80 temporal dynamics of the episode, and how participants would later recount the episode.

81 Results

82 To characterize the shape of the *Sherlock* episode and participants’ subsequent recounts of its
83 unfolding, we used a topic model (?) to discover the latent themes in the episode’s dynamic
84 content. Topic models take as inputs a vocabulary of words to consider and a collection of text
85 documents, and return two output matrices. The first of these is a *topics matrix* whose rows are
86 topics (latent themes) and whose columns correspond to words in the vocabulary. The entries of
87 the topics matrix define how each word in the vocabulary is weighted by each discovered topic.
88 For example, a detective-themed topic might weight heavily on words like “crime,” and “search.”
89 The second output is a *topic proportions matrix*, with one row per document and one column per
90 topic. The topic proportions matrix describes what mixture of discovered topics is reflected in each
91 document.

92 ? collected hand-annotated information about each of 1000 (manually identified) scenes span-
93 ning the roughly 50 minute video used in their experiment. This information included: a brief
94 narrative description of what was happening; whether the scene took place indoors or outdoors;
95 the names of any characters on the screen; the names of any characters who were in focus in the
96 camera shot; the names of characters who were speaking; the location where the scene took place;
97 the camera angle (close up, medium, long, etc.); whether or not background music was present;

98 and other similar details (for a full list of annotated features see *Methods*). We took from these
99 annotations the union of all unique words (excluding stop words, such as “and,” “or,” “but,” etc.)
100 across all features and scenes as the “vocabulary” for the topic model. We then concatenated the
101 sets of words across all features contained in overlapping, 50-scene sliding windows, and treated
102 each 50-scene sequence as a single “document” for the purpose of fitting the topic model. Next,
103 we fit a topic model with (up to) $K = 100$ topics to this collection of documents. We found that
104 32 unique topics (with non-zero weights) were sufficient to describe the time-varying content of
105 the video (see *Methods*; Figs. 1, S2). Note that our approach is similar in some respects to Dy-
106 namic Topic Models (?) in that we sought to characterize how the thematic content of the episode
107 evolved over time. However, whereas Dynamic Topic Models are designed to characterize how
108 the properties of *collections* of documents change over time, our sliding window approach allows
109 us to examine the topic dynamics within a single document (or video). Specifically, our approach
110 yielded (via the topic proportions matrix) a single *topic vector* for each timepoint of the episode (we
111 set timepoints to match the acquisition times of the 1976 fMRI volumes collected as participants
112 viewed the episode).

113 The topics we found were heavily character-focused (e.g., the top-weighted word in each topic
114 was nearly always a character) and could be roughly divided into themes that were primarily
115 Sherlock Holmes-focused (Sherlock is the titular character), primarily John Watson-focused (John
116 is Sherlock’s close confidant and assistant), or focused on Sherlock and John interacting (Fig. S2).
117 Several of the topics were highly similar, which we hypothesized might allow us to distinguish
118 between subtle narrative differences (if the distinctions between those overlapping topics were
119 meaningful; also see Fig. S3). The topic vectors for each timepoint were *sparse*, in that only a small
120 number (usually one or two) of topics tended to be “active” in any given timepoint (Fig. 2A).
121 Further, the dynamics of the topic activations appeared to exhibit *persistiance* (i.e., given that a
122 topic was active in one timepoint, it was likely to be active in the following timepoint) along with
123 *occasional rapid changes* (i.e., occasionally topics would appear to spring into or out of existence).
124 These two properties of the topic dynamics may be seen in the block diagonal structure of the
125 timepoint-by-timepoint correlation matrix (Fig. 2B) and reflect the gradual drift and sudden shifts



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.

fundamental to the contextual dynamics of real-world experiences. Given this observation, we adapted an approach devised by ?, and used a Hidden Markov Model (HMM) to identify the *event boundaries* where the topic activations changed rapidly (i.e., at the boundaries of the blocks in the correlation matrix; event boundaries identified by the HMM are outlined in yellow). Part of our model fitting procedure required selecting an appropriate number of “events” to segment the timeseries into. We used an optimization procedure to identify the number of events that maximized within-event stability while also minimizing across-event correlations (see *Methods* for additional details). To create a stable “summary” of the video, we computed the average topic vector within each event (Fig. 2C).

Given that the time-varying content of the video could be segmented cleanly into discrete events, we wondered whether participants’ recalls of the video also displayed a similar structure. We applied the same topic model (already trained on the video annotations) to each participant’s recalls. Analogous to how we analyzed the time-varying content of the video, to obtain similar estimates for participants’ recalls, we treated each (overlapping) 10-sentence “window” of their transcript as a “document” and then computed the most probable mix of topics reflected in each timepoint’s sentences. This yielded, for each participant, a number-of-windows by number-of-topics topic proportions matrix that characterized how the topics identified in the original video were reflected in the participant’s recalls. Note that an important feature of our approach is that it allows us to compare participant’s recalls to events from the original video, despite that different participants may have used different language to describe the same event, and that those descriptions may not match the original annotations. This is a substantial benefit of projecting the video and recalls into a shared “topic” space. An example topic proportions matrix from one participant’s recalls is shown in Figure 2D.

Although the example participant’s recall topic proportions matrix has some visual similarity to the video topic proportions matrix, the time-varying topic proportions for the example participant’s recalls are not as sparse as for the video (e.g., compare Figs. 2A and D). Similarly, although there do appear to be periods of stability in the recall topic dynamics (e.g., most topics are active or inactive over contiguous blocks of time), the overall timecourses are not as cleanly delineated as

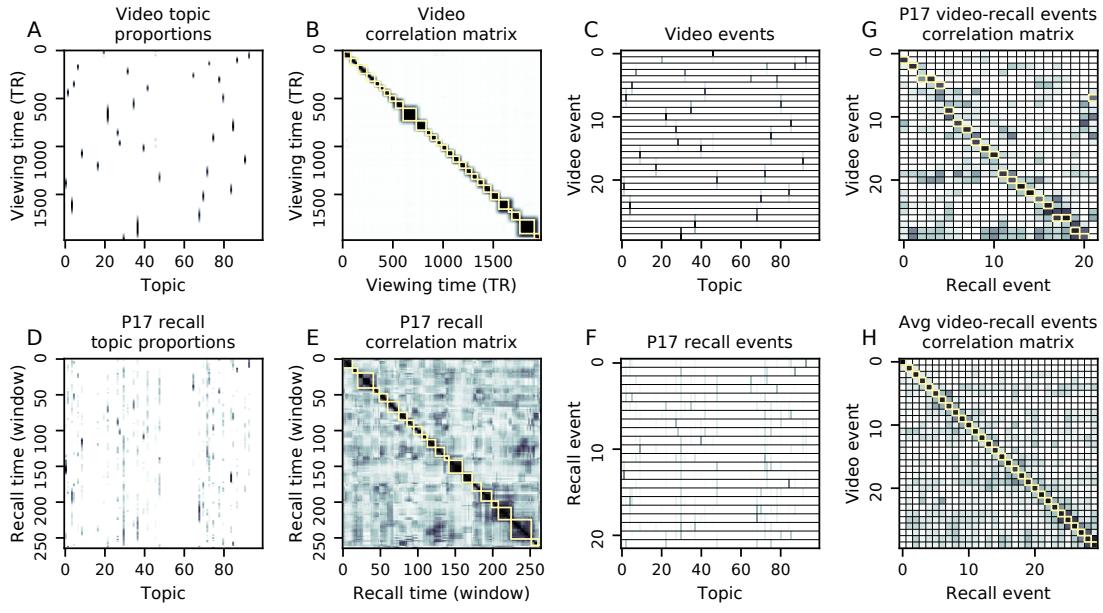


Figure 2: Modelling naturalistic stimuli and recalls. All panels: darker colors indicate greater values; range: [0, 1]. **A.** Topic vectors ($K = 100$) for each of the 1976 video timepoints. **B.** Timepoint-by-timepoint correlation matrix of the topic vectors displayed in Panel A. Event boundaries detected by the HMM are denoted in yellow (30 events detected). **C.** Average topic vectors for each of the 30 video events. **D.** Topic vectors for each of 265 sliding windows of 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 (22 events detected). **F.** Average topic vectors for each of the 22 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 (?).

154 the video topics are. To examine these patterns in detail, we computed the timepoint-by-timepoint
155 correlation matrix for the example participant's recall topic proportions (Fig. 2E). As in the video
156 correlation matrix (Fig. 2B), the example participant's recall correlation matrix has a strong block
157 diagonal structure, indicating that their recalls are discretized into separated events. As for the
158 video correlation matrix, we can use an HMM, along with the aforementioned number-of-events
159 optimization procedure (also see *Methods*) to determine how many events are reflected in the
160 participant's recalls and where specifically the event boundaries fall (outlined in yellow). We
161 carried out a similar analysis on all 17 participants' recall topic proportions matrices (Fig. S4).

162 Two clear patterns emerged from this set of analyses. First, although every individual partic-
163 ipant's recalls could be segmented into discrete events (i.e., every individual participant's recall
164 correlation matrix exhibited clear block diagonal structure; Fig. S4), each participant appeared to
165 have a unique *recall resolution*, reflected in the sizes of those blocks. For example, some participants'
166 recall topic proportions segmented into just a few events (e.g., Participants P4, P5, and P7), while
167 others' recalls segmented into many shorter duration events (e.g., Participants P12, P13, and P17).
168 This suggests that different participants may be recalling the video with different levels of detail-
169 e.g., some might touch on just the major plot points, whereas others might attempt to recall every
170 minor scene or action. The second clear pattern present in every individual participant's recall
171 correlation matrix is that, unlike in the video correlation matrix, there are substantial off-diagonal
172 correlations. Whereas each event in the original video was (largely) separable from the others
173 (Fig. 2B), in transforming those separable events into memory, participants appear to be integrat-
174 ing across multiple events, blending elements of previously recalled and not-yet-recalled events
175 into each newly recalled event (Figs. 2D, S4; also see ??).

176 The above results indicate that both the structure of the original video and participants' recalls
177 of the video exhibit event boundaries that can be identified automatically by characterizing the
178 dynamic content using a shared topic model and segmenting the content into events using HMMs.
179 Next, we asked whether some correspondence might be made between the specific content of the
180 events the participants experienced in the video, and the events they later recalled. One approach
181 to linking the experienced (video) and recalled events is to label each recalled event as matching

182 the video event with the most similar (i.e., most highly correlated) topic vector (Figs. 2G, S5). This
183 yields a sequence of “presented” events from the original video, and a (potentially differently
184 ordered) sequence of “recalled” events for each participant. Analogous to classic list-learning
185 studies, we can then examine participants’ recall sequences by asking which events they tended to
186 recall first (probability of first recall; Fig. 3A; ???); how participants most often transition between
187 recalls of the events as a function of the temporal distance between them (lag-conditional response
188 probability; Fig. 3B; ?); and which events they were likely to remember overall (serial position
189 recall analyses; Fig. 3C; ?). Interestingly, for two of these analyses (probability of first recall and
190 lag-conditional response probability curves) we observe patterns comparable to classic effects from
191 the list-learning literature: namely, a higher probability of initiating recall with the first event in
192 the sequence (Fig. 3A) and a higher probability of transitioning to neighboring events with an
193 asymmetric forward bias (Fig. 3C). In contrast, we do not observe a pattern comparable to the
194 serial position effect (Fig. 3C), but rather we see higher memory for specific events distributed
195 somewhat evenly throughout the video.

196 We can also apply two list-learning-native analyses that describe how participants group items
197 in their recall sequences: temporal clustering and semantic clustering (?; see *Methods* for details).
198 Temporal clustering refers to the extent to which participants group their recall responses according
199 to encoding position. Overall, we found that sequentially viewed video events were clustered
200 heavily in participants’ recall event sequences (mean: 0.767, SEM: 0.029), and that participants
201 with higher temporal clustering scores tended to perform better according to both ?’s hand-
202 annotated memory scores (Pearson’s $r(15) = 0.62$, $p = 0.008$) and our model’s estimate (Pearson’s
203 $r(15) = 0.54$, $p = 0.024$). Semantic clustering measures the extent to which participants cluster
204 their recall responses according to semantic similarity. We found that participants tended to
205 recall semantically similar video events together (mean: 0.787, SEM: 0.018), and that semantic
206 clustering score was also related to both hand-annotated (Pearson’s $r(15) = 0.65$, $p = 0.004$) and
207 model-derived (Pearson’s $r(15) = 0.63$, $p = 0.007$) memory performance.

208 Statistical models of memory studies often treat memory recalls as binary (e.g. the item was re-
209 called or not) and independent events. However, our framework produces a content-based model

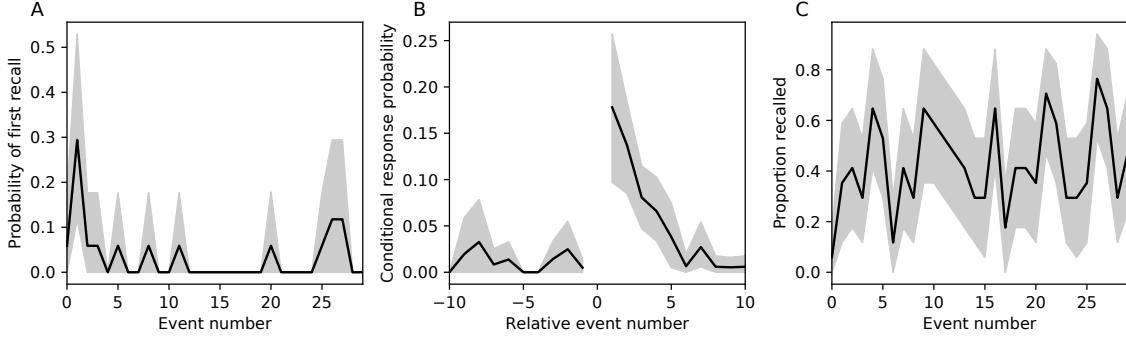


Figure 3: Naturalistic extensions of classic list-learning memory analyses. A. The probability of first recall as a function of the serial position of the event in the video. B. The probability of recalling each event, conditioned on having most recently recalled the event *lag* events away in the video. 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.

of individual stimulus and recall events, allowing for direct quantitative comparison between all stimulus and recall events, as well as between the recall events themselves. Leveraging these content-based models of the stimulus/recall events, we developed two novel metrics for quantifying naturalistic memory representations: *precision* and *distinctiveness*. We define precision as the average correlation between the topic proportions of each recall event and the maximally correlated video event (Fig. 4). Participants whose recall events are more veridical descriptions of what happened in the video event will presumably have higher precision scores. We find that, across participants, a higher precision score is correlated to both hand-annotated memory performance (Pearson's $r(15) = 0.56, p = 0.021$) and the number of recall events estimated by our model (Pearson's $r(15) = 0.85, p < 0.001$). A second novel metric we introduce here is distinctiveness, or how unique the recall description was to each video event. We define distinctiveness as 1 minus the average of all non-matching recall events from the video-recall correlation matrix. We hypothesized that participants who recounted events in a more distinctive way would display better overall memory. We find that this distinctiveness score is related to our model's estimated number of recalled events (Pearson's $r(15) = 0.49, p = 0.046$) but not to the analogous hand-annotated metric (Pearson's $r(15) = 0.31, p = 0.23$). In summary, using two novel metrics afforded by our approach, we find that participants whose recalls are both more precise and distinct remember more content.

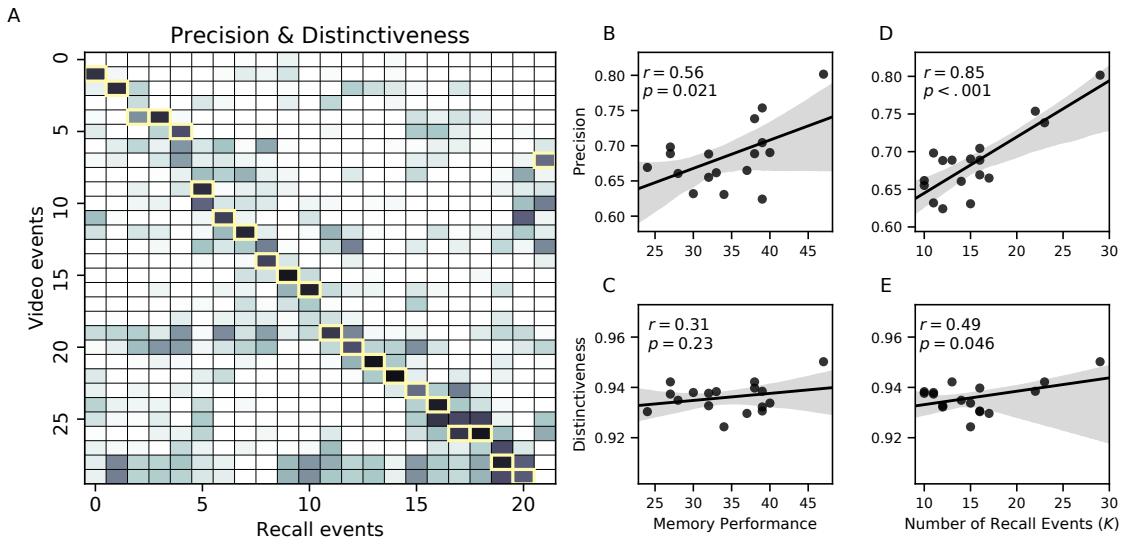


Figure 4: Novel content-based metrics of naturalistic memory: precision and distinctiveness. **A.** A video-recall correlation matrix for a representative participant (17). The yellow boxes highlight the maximum correlation in each column. Precision was computed as the average of the maximum correlation 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).

227 The prior analyses leverage the correspondence between the 100-dimensional topic proportion
228 matrices for the video and participants' recalls to characterize recall. However, it is difficult to gain
229 deep insights into that content solely by examining the topic proportion matrices (e.g., Figs. 2A,
230 D) or the corresponding correlation matrices (Figs. 2B, E, S4). To visualize the time-varying high-
231 dimensional content in a more intuitive way (?) we projected the topic proportions matrices onto a
232 two-dimensional space using Uniform Manifold Approximation and Projection (UMAP; ?). In this
233 lower-dimensional space, each point represents a single video or recall event, and the distances
234 between the points reflect the distances between the events' associated topic vectors (Fig. 5). In
235 other words, events that are near to each other in this space are more semantically similar.

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

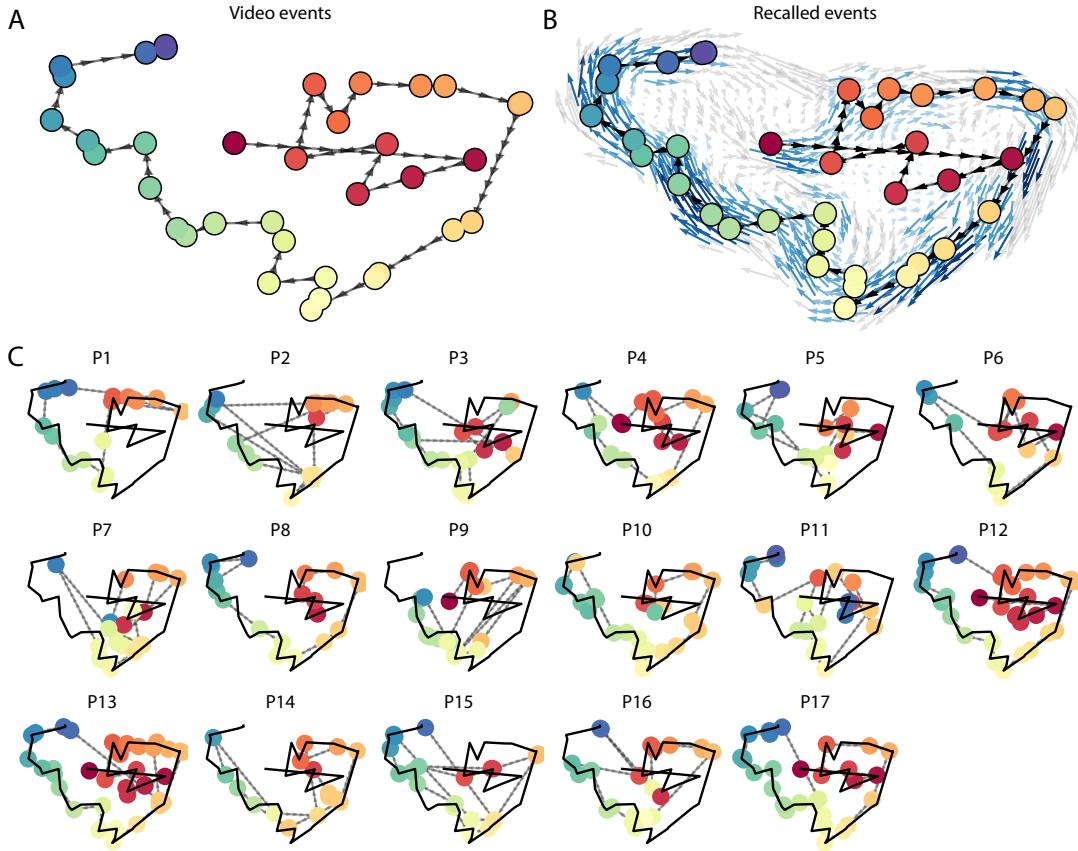


Figure 5: Trajectories through topic space capture the dynamic content of the video and recalls. All panels: the topic proportion matrices have been projected onto a shared two-dimensional space using UMAP. **A.** The two-dimensional topic trajectory taken by the episode of *Sherlock*. Each dot indicates 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.)

255 referenced in participants' recalls) from participants' abilities to recapitulate the full shape of the
256 original video (i.e., the similarity in the shape of the original video trajectory and that defined by
257 each participant's recounting of the video).

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

281 A similar result emerged from assessing the topic vectors for individual video and recall events
282 (Fig. 6C). Here, for each of the three best- and worst-remembered video events, we have constructed

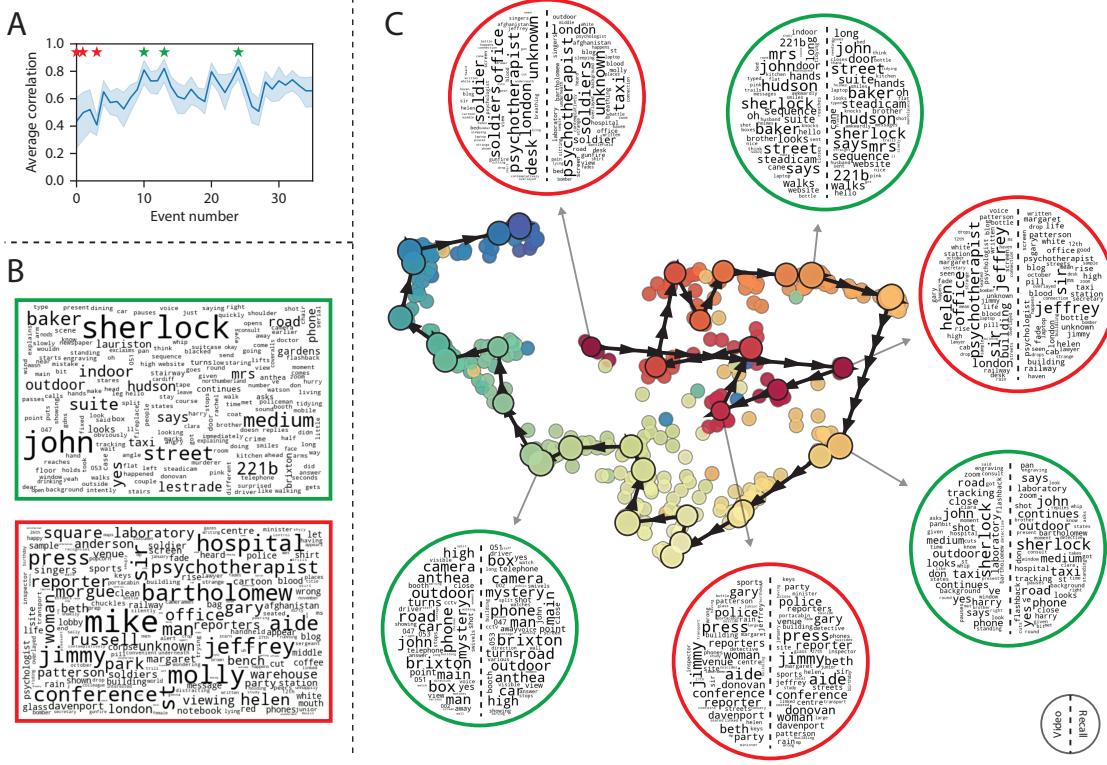


Figure 6: Transforming experience into memory. **A.** Average correlations (across participants) between the topic vectors from each video event and the closest-matching recall events. Error bars denote bootstrap-derived across-participant 95% confidence intervals. The stars denote the three best-remembered 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.

283 two wordles: one from the original video event's topic vector (left) and a second from the average
284 recall topic vector for that event (right). The three best-remembered events (circled in green)
285 correspond to scenes important to the central plot-line (Sherlock and John meeting, Sherlock and
286 John chasing the killer, and the killer calling spying on John in a phone booth). Meanwhile, the three
287 worst-remembered events (circled in red) reflect various side-stories (John talking to his therapist
288 about the war, a soon-to-be-victim's affair with his assistant, and another soon-to-be-victim leaving
289 a party with his friend) that are not essentially to summarizing the video's narrative.

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

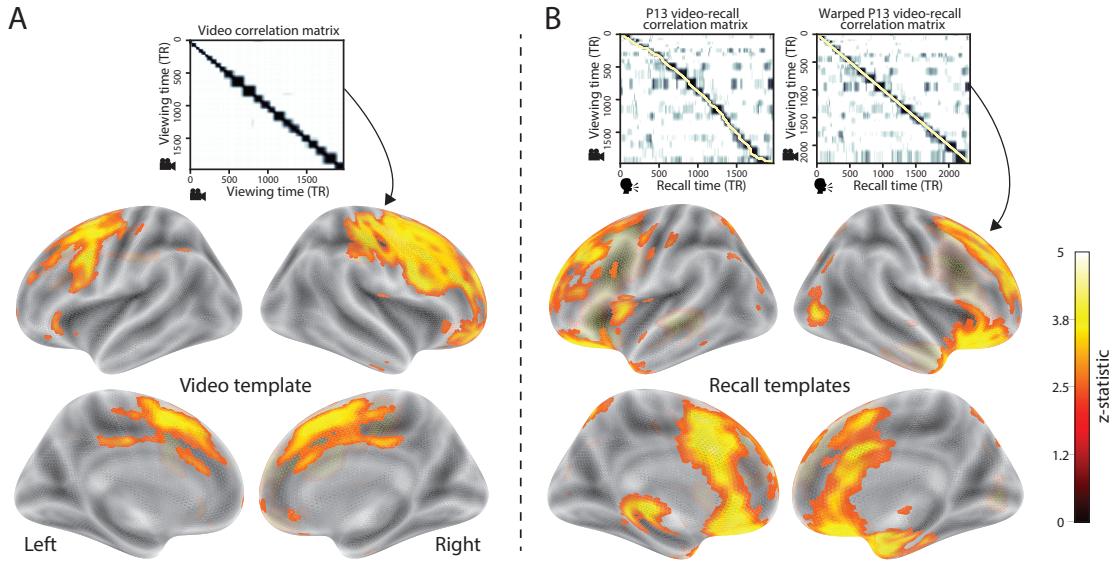


Figure 7: Brain structures that underlie the transformation of experience into memory. **A.** We searched for regions whose responses (as participants watched the video) matched the temporal correlation matrix of the video topic proportions. These regions are sensitive to the narrative structure of the 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.

310 **Discussion**

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

337 person-specific transformations mediated (or represented) by the rMTL and vmPFC may reflect
338 schema knowledge being leveraged, formed, or updated, incorporating ongoing experience into
339 previously acquired knowledge.

340 In extending classical free recall analyses to our naturalistic memory framework, we recovered
341 two patterns of recall dynamics central to list-learning studies: a high probability of initiating
342 recall with the first video event (Fig. 3A) and a strong bias toward transitioning from recalling a
343 given event to recalling the event immediately following it (Fig. 3B). However, equally noteworthy
344 are the typical free recall results not recovered in these analyses, as each highlights a fundamental
345 difference between list-learning studies and naturalistic memory paradigms like the one employed
346 in the present study. The most noticeable departure from hallmark free recall dynamics in these
347 findings is the apparent lack of a serial position effect in Figure 3C, which instead shows greater
348 and lesser recall probabilities for events distributed across the video stimulus. Stimuli in free
349 recall experiments most often comprise lists of simple, common words, presented to participants
350 in a random order. (In fact, numerous word pools have been developed based on these criteria;
351 e.g., ?). These stimulus qualities enable two assumptions that are central to word list analyses,
352 but frequently do not hold for real-world experiences. First, researchers conducting free recall
353 studies may assume that the content at each presentation index is essentially equal, and does not
354 bear qualities that would cause participants to remember it more or less successfully than others.
355 Such is rarely the case with real-world experiences or experiments meant to approximate them,
356 and the effects of both intrinsic and observer-dependent factors on stimulus memorability are well
357 established (for review see ???). Second, the random ordering of list items ensures that (across
358 participants, on average) there is no relationship between the thematic similarity of individual
359 stimuli and their presentation positions—in other words, two semantically related words are no
360 more likely to be presented next to each other than at opposite ends of the list. In most cases, the
361 exact opposite is true of real-world episodes. Our internal thoughts, our actions, and the physical
362 state of the world around us all tend to follow a direct, causal progression. As a result, each moment
363 of our experience tends to be inherently more similar to surrounding moments than to those in
364 the distant past or future. Memory literature has termed this strong temporal autocorrelation

365 “context,” and in various media that depict real-world events (e.g., movies and written stories),
366 we recognize it as a *narrative structure*. While a random word list (by definition) has no such
367 structure, the logical progression between ideas and actions in a naturalistic stimulus prompts the
368 rememberer to recount presented events in order, starting with the beginning. This tendency is
369 reflected in our findings’ second departure from typical free recall dynamics: a lack of increased
370 probability of first recall for end-of-sequence events (Fig. 3A).

371 Thus, analyses such as those in Figure 3 that address only the temporal dynamics of free recall
372 paint an incomplete picture of memory for naturalistic episodes. While useful for studying pre-
373 sentation order-dependent recall dynamics, they neglect to consider the stimuli’s content (or, for
374 example, that content’s potential interrelatedness). However, sensitivity to stimulus and recall con-
375 tent introduces a new challenge: distinguishing between levels of recall quality for a stimulus (i.e.,
376 an event) that is considered to have been “remembered.” When modeling memory experiments,
377 often times events (or items) and their later memories are treated as binary and independent events
378 (e.g., a given list item was simply either remembered or not remembered). Various models of mem-
379 ory (e.g., ?) attempt to improve upon this by including confidence ratings, rendering this binary
380 judgement instead categorical. Our novel framework allows one to assess memory performance in
381 a more continuous way (*precision*), as well as analyze the correlational structure of each encoding
382 event to each memory event (*distinctiveness*). Further and importantly, these two novel metrics we
383 introduce here arise from comparisons of the actual content of the experience/memories, which is
384 not typically modeled. Leveraging this, we find that the successful memory performance is related
385 to 1) the precision with which the participant recounts each event and 2) the distinctiveness of each
386 recall event (relative to the other recalled events). The first finding suggests that the information
387 retained for *any individual event* may predict the overall amount of information retained by the
388 participant. The second finding suggests that the ability to distinguish between temporally or
389 semantically similar content is also related to the quantity of content recovered. Intriguingly, prior
390 studies show that pattern separation, or the ability to discriminate between similar experiences, is
391 impaired in many cognitive disorders as well as natural aging (???). Future work might explore
392 whether and how these metrics compare between cognitively impoverished groups and healthy

393 controls.

394 While a large number of language models exist (e.g., WAS, LSA, word2vec, universal sentence
395 encoder; ???), here we use latent dirichlet allocation (LDA)-based topic models for a few reasons.
396 First, topic models capture the *essence* of a text passage devoid of the specific set and order of words
397 used. This was an important feature of our model since different people may accurately recall a
398 scene using very different language. Second, words can mean different things in different contexts
399 (e.g. “bat” as the act of hitting a baseball, the object used for that action, or as a flying mammal).
400 Topic models are robust to this, allowing words to exist as part of multiple topics. Last, topic models
401 provide a straightforward means to recover the weights for the particular words comprising a topic,
402 enabling easy interpretation of an event’s contents (e.g. Fig. 6). Other models such as Google’s
403 universal sentence encoder offer a context-sensitive encoding of text passages, but the encoding
404 space is complex and non-linear, and thus recovering the original words used to fit the model is
405 not straightforward. However, it’s worth pointing out that our framework is divorced from the
406 particular choice of language model. Moreover, many of the aspects of our framework could be
407 swapped out for other choices. For example, the language model, the timeseries segmentation
408 model and the video-recall matching function could all be customized for the particular problem.
409 Indeed for some problems, recovery of the particular recall words may not be necessary, and thus
410 other text-modeling approaches (such as universal sentence encoder) may be preferable. Future
411 work will explore the influence of particular model choices on the framework’s accuracy.

412 Our work has broad implications for how we characterize and assess memory in real-world
413 settings, such as the classroom or physician’s office. For example, the most commonly used
414 classroom evaluation tools involve simply computing the proportion of correctly answered exam
415 questions. Our work indicates that this approach is only loosely related to what educators might
416 really want to measure: how well did the students understand the key ideas presented in the
417 course? Under this typical framework of assessment, the same exam score of 50% could be
418 ascribed to two very different students: one who attended the full course but struggled to learn
419 more than a broad overview of the material, and one who attended only half of the course but
420 understood the material perfectly. Instead, one could apply our computational framework to build

421 explicit content models of the course material and exam questions. This approach would provide
422 a more nuanced and specific view into which aspects of the material students had learned well
423 (or poorly). In clinical settings, memory measures that incorporate such explicit content models
424 might also provide more direct evaluations of patients' memories.

425 **Methods**

426 **Experimental design and data collection**

427 Data were collected by ?. In brief, participants ($n = 17$) viewed the first 48 minutes of "A Study
428 in Pink", the first episode of the BBC television series *Sherlock*, while fMRI volumes were collected
429 (TR = 1500 ms). The stimulus was divided into a 23 min (946 TR) and a 25 min (1030 TR) segment to
430 mitigate technical issues related to the scanner. After finishing the clip, participants were instructed
431 to (quoting from ?) "describe what they recalled of the [episode] in as much detail as they could, to
432 try to recount events in the original order they were viewed in, and to speak for at least 10 minutes
433 if possible but that longer was better. They were told that completeness and detail were more
434 important than temporal order, and that if at any point they realized they had missed something,
435 to return to it. Participants were then allowed to speak for as long as they wished, and verbally
436 indicated when they were finished (e.g., 'I'm done')." For additional details about the experimental
437 procedure and scanning parameters, see ?. The experimental protocol was approved by Princeton
438 University's Institutional Review Board.

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

445 **Data and code availability**

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

448 **Statistics**

449 All statistical tests we performed were two-sided.

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

451 **Topic modeling**

452 The input to the topic model we trained to characterize the dynamic content of the video comprised
453 hand-generated annotations of each of 1000 scenes spanning the video clip (generated by ?). The
454 features annotated included: narrative details (a sentence or two describing what happened in that
455 scene); whether the scene took place indoors or outdoors; names of any characters that appeared
456 in the scene; name(s) of characters in camera focus; name(s) of characters who were speaking in
457 the scene; the location (in the story) that the scene took place; camera angle (close up, medium,
458 long, top, tracking, over the shoulder, etc.); whether music was playing in the scene or not; and
459 a transcription of any on-screen text. We concatenated the text for all of these features within
460 each segment, creating a “bag of words” describing each scene. We then re-organized the text
461 descriptions into overlapping sliding windows spanning 50 scenes each. In other words, the first
462 text sample comprised the combined text from the first 50 scenes (i.e., 1–50), the second comprised
463 the text from scenes 2–51, and so on. We trained our model using these overlapping text samples
464 with `scikit-learn` (version 0.19.1; ?), called from our high-dimensional visualization and text
465 analysis software, `HyperTools` (?). Specifically, we used the `CountVectorizer` class to transform
466 the text from each scene into a vector of word counts (using the union of all words across all scenes
467 as the “vocabulary,” excluding English stop words); this yielded a number-of-scenes by number-
468 of-words *word count* matrix. We then used the `LatentDirichletAllocation` class (`topics=100`,
469 `method='batch'`) to fit a topic model (?) to the word count matrix, yielding a number-of-scenes

470 (1000) by number-of-topics (100) *topic proportions* matrix. The topic proportions matrix describes
471 which mix of topics (latent themes) is present in each scene. Next, we transformed the topic
472 proportions matrix to match the 1976 fMRI volume acquisition times. For each fMRI volume,
473 we took the topic proportions from whatever scene was displayed for most of that volume's
474 1500 ms acquisition time. This yielded a new number-of-TRs (1976) by number-of-topics (100)
475 topic proportions matrix.

476 We created similar topic proportions matrices using hand-annotated transcripts of each partici-
477 pant's recall of the video (annotated by ?). We tokenized the transcript into a list of sentences, and
478 then re-organized the list into overlapping sliding windows spanning 10 sentences each; in turn
479 we transformed each window's sentences into a word count vector (using the same vocabulary as
480 for the video model). We then used the topic model already trained on the video scenes to compute
481 the most probable topic proportions for each sliding window. This yielded a number-of-sentences
482 (range: 68–294) by number-of-topics (100) topic proportions matrix, for each participant. These
483 reflected the dynamic content of each participant's recalls. Finally, we resampled each recall model
484 to match the timecourse of the video model. Note: for details on how we selected the video and
485 recall window lengths and number of topics, see *Supporting Information* and Figure S1.

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

487 We parsed the topic trajectories of the video and participants' recalls into events using Hidden
488 Markov Models (?). Given the topic proportions matrix (describing the mix of topics at each
489 timepoint) and a number of states, K , an HMM recovers the set of state transitions that segments
490 the timeseries into K discrete states. Following ?, we imposed an additional set of constraints on
491 the discovered state transitions that ensured that each state was encountered exactly once (i.e.,
492 never repeated). We used the BrainIAK toolbox (?) to implement this segmentation.

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

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

495 where a was the average correlation between the topic vectors of timepoints within the same state;
496 b was the average correlation between the topic vectors of timepoints within *different* states; J was
497 a constant used to ensure a positive denominator, set equal to $\min_K \left[\frac{a}{b} \right]$; and α was a regularization
498 parameter that we set to 5 times the window length (i.e., 250 scenes for the video topic trajectory
499 and 50 sentences for the recall topic trajectories). Before subtracting the regularization term, we
500 scaled the ratio of correlations to be a proportion of the maximum ratio across all K 's. Figure 2B
501 displays the event boundaries returned for the video, and Figure S4 displays the event boundaries
502 returned for each participant's recalls (See Fig. S6 for the values of a and b for each K , Fig. S7 for
503 the optimization functions for the video and recalls). After obtaining these event boundaries, we
504 created stable estimates of each topic proportions matrix by averaging the topic vectors within
505 each event. This yielded a number-of-events by number-of-topics matrix for the video and recalls
506 from each participant.

507 We also evaluated a parameter-free procedure for choosing K , which finds the K value that
508 maximizes the Wasserstein distance (a.k.a. “Earth mover’s” distance) between the within and
509 across event distributions of correlation values. This alternative procedure largely replicated the
510 pattern of results found with the parameterized method described above, but recovered sub-
511 stantially fewer events on average (Fig.S8). While both approaches seem to underestimate the
512 number of video/recall events relative to the “true” number (as determined by human raters), the
513 parameterized approach was closer to the true number.

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

515 In traditional list-learning experiments, participants view a list of items (e.g., words) and then recall
516 the items later. Our video-recall event matching approach affords us the ability to analyze memory
517 in a similar way. The video and recall events can be treated analogously to studied and recalled
518 “items” in a list-learning study. We can then extend classic analyses of memory performance and
519 dynamics (originally designed for list-learning experiments) to the more naturalistic video recall
520 task used in this study.

521 Perhaps the simplest and most widely used measure of memory performance is *accuracy*—i.e.,

522 the proportion of studied (experienced) items (in this case, the 34 video events) that the participant
523 later remembered. ? developed a human rating system whereby the quality of each participant's
524 memory was evaluated by an independent rater. We found a strong across-participants correlation
525 between these independant ratings and the overall number of events that our HMM approach
526 identified in participants' recalls (Pearson's $r(15) = 0.64, p = 0.006$).

527 As described below, we next considered a number of memory performance measures that are
528 typically associated with list-learning studies. We also provide a software package, Quail, for
529 carrying out these analyses (?).

530 **Probability of first recall (PFR).** PFR curves (???) reflect the probability that an item will be
531 recalled first as a function of its serial position during encoding. To carry out this analysis, we
532 initialized a number-of-participants (17) by number-of-video-events (34) matrix of zeros. Then for
533 each participant, we found the index of the video event that was recalled first (i.e., the video event
534 whose topic vector was most strongly correlated with that of the first recall event) and filled in that
535 index in the matrix with a 1. Finally, we averaged over the rows of the matrix, resulting in a 1 by
536 34 array representing the proportion of participants that recalled an event first, as a function of the
537 order of the event's appearance in the video (Fig. 3A).

538 **Lag conditional probability curve (lag-CRP).** The lag-CRP curve (?) reflects the probability of
539 recalling a given event after the just-recalled event, as a function of their relative positions (or *lag*).
540 In other words, a lag of 1 indicates that a recalled event came immediately after the previously
541 recalled event in the video, and a lag of -3 indicates that a recalled event came 3 events before the
542 previously recalled event. For each recall transition (following the first recall), we computed the
543 lag between the current recall event and the next recall event, normalizing by the total number of
544 possible transitions. This yielded a number-of-participants (17) by number-of-lags (-33 to +33; 67
545 lags total) matrix. We averaged over the rows of this matrix to obtain a group-averaged lag-CRP
546 curve (Fig. 3B).

547 **Serial position curve (SPC).** SPCs (?) reflect the proportion of participants that remember each
548 item as a function of the items' serial position during encoding. We initialized a number-of-
549 participants (17) by number-of-video-events (34) matrix of zeros. Then, for each recalled event,
550 for each participant, we found the index of the video event that the recalled event most closely
551 matched (via the correlation between the events' topic vectors) and entered a 1 into that position
552 in the matrix (i.e., for the given participant and event). This resulted in a matrix whose entries
553 indicated whether or not each event was recalled by each participant (depending on whether the
554 corresponding entires were set to one or zero). Finally, we averaged over the rows of the matrix
555 to yield a 1 by 34 array representing the proportion of participants that recalled each event as a
556 function of the order of the event's appearance in the video (Fig. 3C).

557 **Temporal clustering scores.** Temporal clustering describes participants' tendency to organize
558 their recall sequences by the learned items' encoding positions. For instance, if a participant
559 recalled the video events in the exact order they occurred (or in exact reverse order), this would
560 yield a score of 1. If a participant recalled the events in random order, this would yield an expected
561 score of 0.5. For each recall event transition (and separately for each participant), we sorted
562 all not-yet-recalled events according to their absolute lag (i.e., distance away in the video). We
563 then computed the percentile rank of the next event the participant recalled. We averaged these
564 percentile ranks across all of the participant's recalls to obtain a single temporal clustering score
565 for the participant.

566 **Semantic clustering scores.** Semantic clustering describes participants' tendency to recall seman-
567 tically similar presented items together in their recall sequences. Here, we used the topic vectors
568 for each event as a proxy for its semantic content. Thus, the similarity between the semantic
569 content for two events can be computed by correlating their respective topic vectors. For each
570 recall event transition, we sorted all not-yet-recalled events according to how correlated the topic
571 vector of the closest-matching video event was to the topic vector of the closest-matching video event
572 to the just-recalled event. We then computed the percentile rank of the observed next recall. We

573 averaged these percentile ranks across all of the participant's recalls to obtain a single semantic
574 clustering score for the participant.

575 **Novel naturalistic memory metrics**

576 **Precision.** We tested whether participants who recalled more events were also more *precise* in
577 their recollections. For each participant, we computed the average correlation between the topic
578 vectors for each recall event and those of its closest-matching video event. This gave a single value
579 per participant representing the average precision across all recalled events. We then Fisher's *z*-
580 transformed these values and correlated them with both hand-annotated and model-derived (i.e.,
581 k or the number of events recovered by the HMM) memory performance.

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

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

588 where $\bar{c}(\text{event})$ is the average correlation between the given recalled event's topic vector and the
589 topic vectors from all video events *except* the best-matching video event. We then averaged these
590 distinctiveness scores across all of the events recalled by the given participant. As above, we used
591 Fisher's *z*-transformation before correlating these values with hand-annotated and model derived
592 memory performance scores across-subjects.

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

594 We used the UMAP algorithm (?) to project the 100-dimensional topic space onto a two-dimensional
595 space for visualization (Figs. 5, 6). To ensure that all of the trajectories were projected onto the *same*
596 lower dimensional space, we computed the low-dimensional embedding on a “stacked” matrix
597 created by vertically concatenating the events-by-topics topic proportions matrices for the video
598 and all 17 participants’ recalls. We then divided the rows of the result (a total-number-of-events
599 by two matrix) back into separate matrices for the video topic trajectory and the trajectories for
600 each participant’s recalls (Fig. 5). This general approach for discovering a shared low-dimensional
601 embedding for a collections of high-dimensional observations follows ?.

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

603 In Figure 5B, we present an analysis aimed at characterizing locations in topic space that dif-
604 ferent participants move through in a consistent way (via their recall topic trajectories). The
605 two-dimensional topic space used in our visualizations (Fig. 5) comprised a 9x9 (arbitrary units)
606 square. We divided this space into a grid of vertices spaced 0.25 units apart. For each vertex, we
607 examined the set of line segments formed by connecting each pair successively recalled events,
608 across all participants, that passed within 0.5 units. We computed the distribution of angles formed
609 by those segments and the x -axis, and used a Rayleigh test to determine whether the distribution
610 of angles was reliably “peaked” (i.e., consistent across all transitions that passed through that local
611 portion of topic space). To create Figure 5B we drew an arrow originating from each grid vertex,
612 pointing in the direction of the average angle formed by line segments that passed within 0.5 units.
613 We set the arrow lengths to be inversely proportional to the p -values of the Rayleigh tests at each
614 vertex. Specifically, for each vertex we converted all of the angles of segments that passed within
615 0.5 units to unit vectors, and we set the arrow lengths at each vertex proportional to the length
616 of the (circular) mean vector. We also indicated any significant results ($p < 0.05$, corrected using
617 the Benjamani-Hochberg procedure) by coloring the arrows in blue (darker blue denotes a lower
618 p -value, i.e., a longer mean vector); all tests with $p \geq 0.05$ are displayed in gray and given a lower

619 opacity value.

620 **Searchlight fMRI analyses**

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

628 Next, we constructed two sets of “template” matrices: one reflecting the video’s topic trajectory
629 and the other reflecting each participant’s recall topic trajectory. To construct the video template, we
630 computed the correlations between the topic proportions estimated for every pair of TRs (prior to
631 segmenting the trajectory into discrete events; i.e., the correlation matrix shown in Figs. 2B and 7A).
632 We constructed similar temporal correlation matrices for each participant’s recall topic trajectory
633 (Figs. 2D, S4). However, to correct for length differences and potential non-linear transformations
634 between viewing time and recall time, we first used dynamic time warping (?) to temporally align
635 participants’ recall topic trajectories with the video topic trajectory (an example correlation matrix
636 before and after warping is shown in Fig. 7B). This yielded a 1976 by 1976 correlation matrix for
637 the video template and for each participant’s recall template.

638 To determine which (cubes of) voxel responses reliably matched the video template, we cor-
639 related the upper triangle of the voxel correlation matrix for each cube with the upper triangle
640 of the video template matrix (?). This yielded, for each participant, a single correlation value.
641 We computed the average (Fisher z-transformed) correlation coefficient across participants. We
642 used a permutation-based procedure to assess significance, whereby we re-computed the average
643 correlations for each of 100 “null” video templates (constructed by circularly shifting the template
644 by a random number of timepoints). (For each permutation, the same shift was used for all partici-
645 pants.) We then estimated a p -value by computing the proportion of shifted correlations that were

646 larger than the observed (unshifted) correlation. To create the map in Figure 7A we thresholded
647 out any voxels whose correlation values fell below the 95th percentile of the permutation-derived
648 null distribution.

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

658 **References**

659 **Supporting information**

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

661 **Acknowledgements**

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

667 **Author contributions**

668 Conceptualization: A.C.H. and J.R.M.; Methodology: A.C.H. and J.R.M.; Software: A.C.H., P.C.F.
669 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.
670 and J.R.M.; Supervision: J.R.M.

671 **Author information**

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