

1 Geometric models reveal behavioral and neural
2 signatures of how naturalistic experiences are
3 transformed into episodic memories

4 Andrew C. Heusser^{1, 2, †}, Paxton C. Fitzpatrick^{1, †}, and Jeremy R. Manning^{1, *}

¹Department of Psychological and Brain Sciences

Dartmouth College, Hanover, NH 03755, USA

²Akili Interactive

Boston, MA 02110

[†]Denotes equal contribution

^{*}Corresponding author: Jeremy.R.Manning@dartmouth.edu

5 August 19, 2020

6 **Abstract**

7 Our ongoing subjective experience reflects external sensory information from each moment,
8 along with additional information from our past that we carry with us into that moment. The
9 blend of memories, knowledge, emotions, goals, and other internal perceptual and mental states
10 that color our subjective experience provides a *context* for interpreting new information and
11 conceptually linking what is happening now with our prior experiences. Because this context-
12 ual information is often person-specific, the subjective experience that each person encodes into
13 their memory is often idiosyncratic, even for shared experiences and sensory perspectives. We
14 sought to study which aspects of a shared naturalistic experience were preserved or distorted,
15 and how those distortions compared across individuals. To this end, we developed a geomet-

ric framework for mathematically characterizing the subjective conceptual content of dynamic naturalistic experiences. We model experiences as *trajectories* through word embedding spaces whose coordinates reflect the universe of thoughts under consideration. We also demonstrate how *memories* may also be modeled as trajectories through the same spaces. According to this view, encoding an experience into memory entails geometrically distorting or transforming the original experience’s trajectory. This translates qualitative neuropsychological questions about how we remember naturalistic experiences into quantitative geometric questions about the spatial configurations of trajectory shapes. We applied our framework to data collected as participants watched and verbally recounted a television episode while undergoing functional neuroimaging. We found that the trajectories of participants’ recounts of the episode nearly all captured the coarse spatial properties of the original episode’s trajectory (i.e., the essential plot points), but participants differed in their memory for fine details. We also identified a network of brain structures that were sensitive to the shape of the episode’s trajectory through word embedding space, and an overlapping network that predicted, at the time of encoding, how people would distort (transform) the episode’s trajectory when they recounted the episode later. Our work provides insights into how our brains distort and transform our ongoing experiences when we encode them into episodic memories.

33 **Introduction**

34 What does it mean to *remember* something? In traditional episodic memory experiments (e.g.,
35 list-learning or trial-based experiments; Murdock, 1962; Kahana, 1996), remembering is often cast
36 as a discrete and binary operation: each studied item may be separated from the rest of one’s
37 experience and singularly labeled as having been recalled or forgotten. More nuanced studies
38 might incorporate self-reported confidence measures as a proxy for memory strength, or ask
39 participants to discriminate between “recollecting” the (contextual) details of an experience or
40 having a general feeling of “familiarity” (Yonelinas, 2002). Using well-controlled, trial-based
41 experimental designs, the field has amassed a wealth of information regarding human episodic
42 memory. However, there are fundamental properties of the external world and our memories that

43 trial-based experiments are not well-suited to capture (for review, also see Koriat and Goldsmith,
44 1994; Huk et al., 2018). First, our experiences and memories are continuous, rather than discrete—
45 isolating a (naturalistic) event from the context in which it occurs can substantially change its
46 meaning. Second, asking whether the rememberer has precisely reproduced a specific set of
47 words to describe a given experience is nearly orthogonal to how well they were actually able to
48 remember it. In classic (e.g., list-learning) memory studies, by contrast, the number or proportion
49 of (exact) recalls is often considered to be a primary metric for assessing the quality of participants'
50 memories. Third, one might remember the *essence* (or a general summary) of an experience but
51 forget (or neglect to recount) particular details. Capturing the essence of what happened is often
52 a main goal of recounting an episodic memory to a listener, whereas the inclusion of specific
53 low-level details is often less pertinent.

54 How might we formally characterize the *essence* of an experience, and whether it has been re-
55 covered by the rememberer? And how might we distinguish an experience's overarching essence
56 from its low-level details? One approach is to start by considering some fundamental proper-
57 ties of the dynamics of our experiences. Any given moment of an experience tends to derive
58 meaning from surrounding moments, as well as from longer-range temporal associations (Lerner
59 et al., 2011; Manning, 2019, 2020). Therefore, the timecourse describing how an event unfolds is
60 fundamental to its overall meaning. Further, this hierarchy formed by our subjective experiences
61 at different timescales defines a *context* for each new moment (e.g., Howard and Kahana, 2002;
62 Howard et al., 2014), and plays an important role in how we interpret that moment and remember
63 it later (for review see Manning et al., 2015; Manning, 2020). Our memory systems can leverage
64 these associations to form predictions that help guide our behaviors (Ranganath and Ritchey,
65 2012). For example, as we navigate the world, the features of our subjective experiences tend
66 to change gradually (e.g., the room or situation we are in at any given moment is strongly tem-
67 porally autocorrelated), allowing us to form stable estimates of our current situation and behave
68 accordingly (Zacks et al., 2007; Zwaan and Radvansky, 1998).

69 Occasionally, this gradual “drift” of our ongoing experience is punctuated by sudden changes,
70 or “shifts” (e.g., when we walk through a doorway; Radvansky and Zacks, 2017). Prior research

71 suggests that these sharp transitions (termed *event boundaries*) help to discretize our experiences
72 (and their mental representations) into *events* (Radvansky and Zacks, 2017; Brunec et al., 2018;
73 Heusser et al., 2018a; Clewett and Davachi, 2017; Ezzyat and Davachi, 2011; DuBrow and Davachi,
74 2013). The interplay between the stable (within-event) and transient (across-event) temporal
75 dynamics of an experience also provides a potential framework for transforming experiences
76 into memories that distills those experiences down to their essence. For example, prior work
77 has shown that event boundaries can influence how we learn sequences of items (Heusser et al.,
78 2018a; DuBrow and Davachi, 2013), navigate (Brunec et al., 2018), and remember and understand
79 narratives (Zwaan and Radvansky, 1998; Ezzyat and Davachi, 2011). This work also suggests a
80 means of distinguishing the essence of an experience from its low-level details. The event-level
81 properties reflect how the high-level experience unfolds (i.e., its essence), whereas subtler within-
82 event changes reflect low-level details. Prior research has also implicated a network of brain
83 regions (including the hippocampus and the medial prefrontal cortex) as playing a critical role
84 in transforming experiences into structured and consolidated memories (Tomparay and Davachi,
85 2017).

86 Here, we sought to examine how the temporal dynamics of a “naturalistic” experience were later
87 reflected in participants’ memories. We also sought to leverage the above conceptual insights into
88 the distinctions between an experience’s essence versus its low-level details to build models that
89 explicitly quantified these distinctions. We analyzed an open dataset that comprised behavioral
90 and functional Magnetic Resonance Imaging (fMRI) data collected as participants viewed and then
91 verbally recounted an episode of the BBC television series *Sherlock* (Chen et al., 2017). We developed
92 a computational framework for characterizing the temporal dynamics of the moment-by-moment
93 content of the episode, and of participants’ verbal recalls. Specifically, we use topic modeling (Blei
94 et al., 2003) to characterize the thematic conceptual (semantic) content present in each moment
95 of the episode and recalls, and hidden Markov models (Rabiner, 1989; Baldassano et al., 2017) to
96 discretize this evolving semantic content into events. In this way, we cast naturalistic experiences
97 (and memories of those experiences) as geometric *trajectories* that describe how the experiences
98 evolve over time. Under this framework, successful remembering entails verbally “traversing”

99 the content trajectory of the episode, thereby reproducing the shape (or essence) of the original
100 experience. The essence of the episode is captured by the sequence of geometric coordinates of its
101 events, and the episode's low-level details are captured by examining its within-event geometric
102 properties.

103 Comparing the overall shapes of the topic trajectories of the episode and of participants'
104 retellings of the episode reveals which aspects of the episode's essence were preserved (or dis-
105 carded) in the translation into memory. We also develop two metrics for assessing participants'
106 memories for low-level details: (1) the *precision* with which a participant recounts details about
107 each event, and (2) the *distinctiveness* of each recall event (relative to other recalled events). We
108 examine how these metrics relate to overall memory performance as judged by third-party human
109 annotators. We also compare and contrast our general approach to studying memory for naturalis-
110 tic experiences with standard metrics for assessing performance on more traditional memory tasks,
111 such as list learning. Last, we leverage our framework to identify networks of brain structures
112 whose responses (as participants watched the episode) reflected the temporal dynamics of either
113 the episode or how participants would later recount it.

114 **Results**

115 To characterize the dynamic content of the *Sherlock* episode and participants' subsequent recount-
116 ings, we used a topic model (Blei et al., 2003) to discover the episode's latent themes. Topic models
117 take as inputs a vocabulary of words to consider and a collection of text documents, and return two
118 output matrices. The first of these is a *topics matrix* whose rows are *topics* (latent themes) and whose
119 columns correspond to words in the vocabulary. The entries of the topics matrix reflect how each
120 word in the vocabulary is weighted by each discovered topic. For example, a detective-themed
121 topic might weight heavily on words like "crime," and "search." The second output is a *topic*
122 *proportions matrix*, with one row per document and one column per topic. The topic proportions
123 matrix describes what mixture of discovered topics is reflected in each document.

124 Chen et al. (2017) collected hand-annotated information about each of 1000 (manually identified)
125 scenes spanning the roughly 50 minute video used in their experiment. This information included:
126 a brief narrative description of what was happening, the location where the scene took place, the
127 names of any characters on the screen, and other similar details (for a full list of annotated features,
128 see *Methods*). We took from these annotations the union of all unique words (excluding stop
129 words, such as “and,” “or,” “but,” etc.) across all features and scenes as the “vocabulary” for the
130 topic model. We then concatenated the sets of words across all features contained in overlapping,
131 sliding windows of (up to) 50 scenes, and treated each window as a single “document” for the
132 purpose of fitting the topic model. Next, we fit a topic model with (up to) $K = 100$ topics to this
133 collection of documents. We found that 32 unique topics (with non-zero weights) were sufficient
134 to describe the time-varying content of the video (see *Methods*; Figs. 1, S2). Note that our approach
135 is similar in some respects to Dynamic Topic Models (Blei and Lafferty, 2006) in that we sought
136 to characterize how the thematic content of the episode evolved over time. However, whereas
137 Dynamic Topic Models are designed to characterize how the properties of *collections* of documents
138 change over time, our sliding window approach allows us to examine the topic dynamics within
139 a single document (or video). Specifically, our approach yielded (via the topic proportions matrix)
140 a single *topic vector* for each sliding window of annotations transformed by the topic model. We
141 then stretched (interpolated) the resulting windows-by-topics matrix to match the time series of
142 the 1976 fMRI volumes collected as participants viewed the episode.

143 The 32 topics we found were heavily character-focused (i.e., the top-weighted word in each
144 topic was nearly always a character) and could be roughly divided into themes centered around
145 Sherlock Holmes (the titular character), John Watson (Sherlock’s close confidant and assistant),
146 supporting characters (e.g., Inspector Lestrade, Sergeant Donovan, or Sherlock’s brother Mycroft),
147 or the interactions between various groupings of these characters (see Fig. S2). Several of the
148 identified topics were highly similar, which we hypothesized might allow us to distinguish between
149 subtle narrative differences if the distinctions between those overlapping topics were meaningful.
150 The topic vectors for each timepoint were also *sparse*, in that only a small number (usually one
151 or two) of topics tended to be “active” in any given timepoint (Fig. 2A). Further, the dynamics

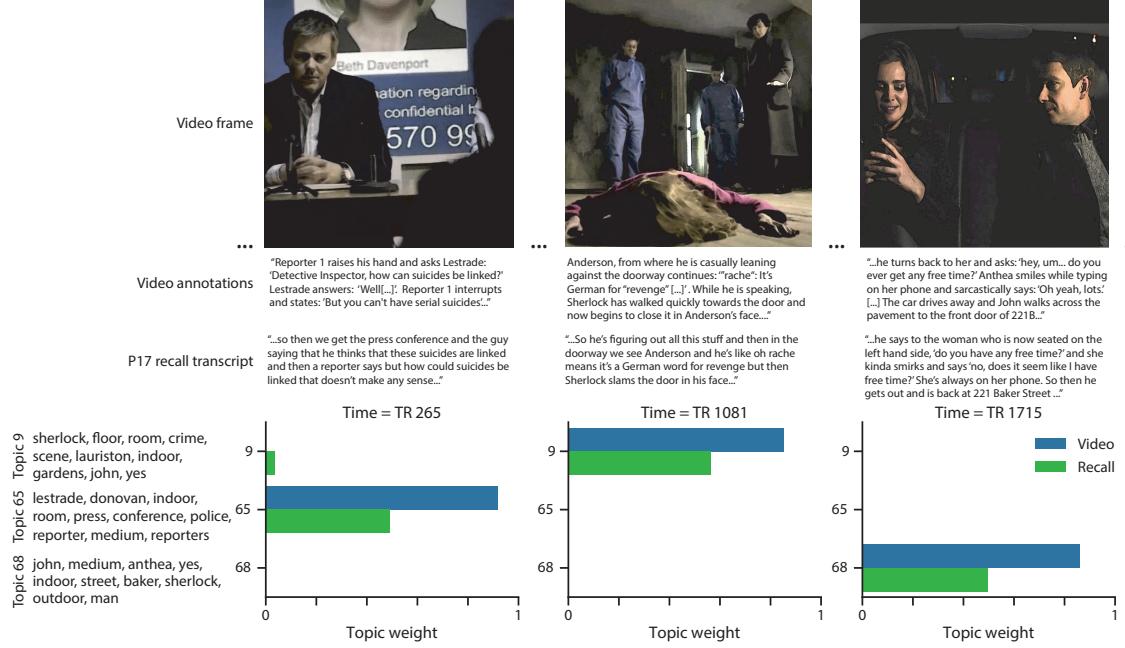


Figure 1: Methods overview. We used hand-annotated descriptions of each moment of video to fit a topic model. Three example video frames and their associated descriptions are displayed (top two rows). Participants later recalled the video (in the third row, we show example recalls of the same three scenes from participant 17). 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 10 highest-weighted words for each topic. Figure S2 provides a full list of the top 10 words from each of the discovered topics.

of the topic activations appeared to exhibit *persistence* (i.e., given that a topic was active in one timepoint, it was likely to be active in the following timepoint) along with *occasional rapid changes* (i.e., occasionally topics would appear to spring into or out of existence). These two properties of the topic dynamics may be seen in the block diagonal structure of the timepoint-by-timepoint correlation matrix (Fig. 2B) and reflect the gradual drift and sudden shifts fundamental to the temporal dynamics of real-world experiences. Given this observation, we adapted an approach devised by Baldassano et al. (2017), and used a hidden Markov model (HMM) to identify the *event boundaries* where the topic activations changed rapidly (i.e., the boundaries of the blocks in the temporal correlation matrix; event boundaries identified by the HMM are outlined in yellow in Fig. 2B). Part of our model fitting procedure required selecting an appropriate number of “events” into which the topic trajectory should be segmented. To accomplish this, we used an optimization procedure that maximized the difference between the topic weights for timepoints within an event versus timepoints across multiple events (see *Methods* for additional details). We then created a stable “summary” of the content within each video event by averaging the topic vectors across the timepoints spanned by 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. Analogously to how we parsed the time-varying content of the video, to obtain similar estimates for each participant’s recall, we treated each overlapping window of (up to 10) sentences from their transcript as a “document,” and 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 participants’ recalls to events from the original video, despite different participants using widely varying language to describe the events, and that those descriptions often diverged in content and quality from the video annotations. This is a substantial benefit of projecting the video and recalls into a shared “topic” space. An example topic proportions matrix

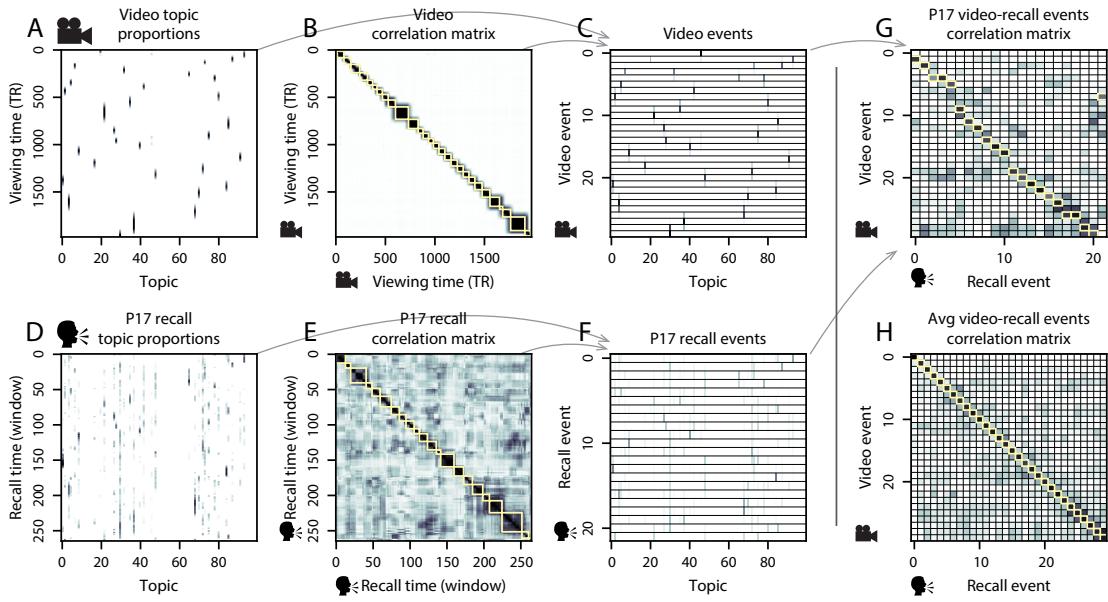


Figure 2: Modelling naturalistic stimuli and recalls. All panels: darker colors indicate greater values; range: [0, 1]. **A.** Topic vectors ($K = 100$) for each of the 1976 video timepoints. **B.** Timepoint-by-timepoint correlation matrix of the topic vectors displayed in Panel A. Event boundaries discovered 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). For similar plots for all participants, see Figure S4. **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).

180 from one participant’s recalls is shown in Figure 2D.

181 Although the example participant’s recall topic proportions matrix has some visual similarity to
182 the video topic proportions matrix, the time-varying topic proportions for the example participant’s
183 recalls are not as sparse as those for the video (compare Figs. 2A and D). Similarly, although
184 there do appear to be periods of stability in the recall topic dynamics (i.e., most topics are active
185 or inactive over contiguous blocks of time), the changes in topic activations that define event
186 boundaries appear less clearly delineated in participants’ recalls than in the episode’s annotations.
187 To examine these patterns in detail, we computed the timepoint-by-timepoint correlation matrix
188 for the example participant’s recall trajectory (Fig. 2E). As in the video correlation matrix (Fig. 2B),
189 the example participant’s recall correlation matrix has a strong block diagonal structure, indicating
190 that their recalls are discretized into separated events. As for the video correlation matrix, we
191 leveraged an HMM-based optimization procedure (see *Methods*) to determine how many events
192 are reflected in the participant’s recalls and where specifically the event boundaries fall (outlined
193 in yellow). We carried out a similar analysis on all 17 participants’ recall topic proportions matrices
194 (Fig. S4).

195 Two clear patterns emerged from this set of analyses. First, although every individual partic-
196 ipant’s recalls could be segmented into discrete events (i.e., every individual participant’s recall
197 correlation matrix exhibited clear block diagonal structure; Fig. S4), each participant appeared to
198 have a unique *recall resolution*, reflected in the sizes of those blocks. While some participants’ recall
199 topic proportions segmented into just a few events (e.g., Participants P4, P5, and P7), others’ seg-
200 mented into many shorter duration events (e.g., Participants P12, P13, and P17). This suggests that
201 different participants may be recalling the video with different levels of detail—i.e., some might
202 touch on just the major plot points, whereas others might attempt to recall every minor scene or
203 action. The second clear pattern present in every individual participant’s recall correlation matrix
204 was that, unlike in the video correlation matrix, there were substantial off-diagonal correlations.
205 Whereas each event in the original video was (largely) separable from the others (Fig. 2B), in
206 transforming those separable events into memory, participants appeared to be integrating across
207 multiple events, blending elements of previously recalled and not-yet-recalled content into each

208 newly recalled event (Figs. 2E, S4; also see Manning et al., 2011; Howard et al., 2012; Manning,
209 2019).

210 The above results indicate that both the structure of the original video and participants' recalls
211 of the video exhibit event boundaries that can be identified automatically by characterizing the
212 dynamic content using a shared topic model and segmenting the content into events via HMMs.
213 Next, we asked whether some correspondence might be made between the specific content of the
214 events the participants experienced in the video, and the events they later recalled. One approach
215 to linking the experienced (video) and recalled events is to label each recalled event as matching
216 the video event with the most similar (i.e., most highly correlated) topic vector (Figs. 2G, S5). This
217 yields a sequence of "presented" events from the original video, and a (potentially differently
218 ordered) sequence of "recalled" events for each participant. Analogous to classic list-learning
219 studies, we can then examine participants' recall sequences by asking which events they tended
220 to recall first (probability of first recall; Fig. 3A; Atkinson and Shiffrin, 1968; Postman and Phillips,
221 1965; Welch and Burnett, 1924); how participants most often transition between recalls of the
222 events as a function of the temporal distance between them (lag-conditional response probability;
223 Fig. 3B; Kahana, 1996); and which events they were likely to remember overall (serial position
224 recall analyses; Fig. 3C; Murdock, 1962). Interestingly, for two of these analyses (probability of
225 first recall and lag-conditional response probability curves) we observed patterns comparable to
226 classic effects from list-learning literature: namely, a higher probability of initiating recall with the
227 first event in the sequence (Fig. 3A) and a higher probability of transitioning to neighboring events
228 with an asymmetric forward bias (Fig. 3B). In contrast, we did not observe a pattern comparable
229 to the serial position effect (Fig. 3C), but rather greater memory for specific events distributed
230 approximately evenly throughout the video.

231 We can also apply two list-learning-native analyses that describe how participants group items
232 in their recall sequences: temporal clustering and semantic clustering (Polyn et al., 2009, see
233 *Methods* for details). Temporal clustering refers to the extent to which participants group their
234 recall responses according to encoding position. Overall, we found that sequentially viewed video
235 events were clustered heavily in participants' recall event sequences (mean clustering score: 0.767,

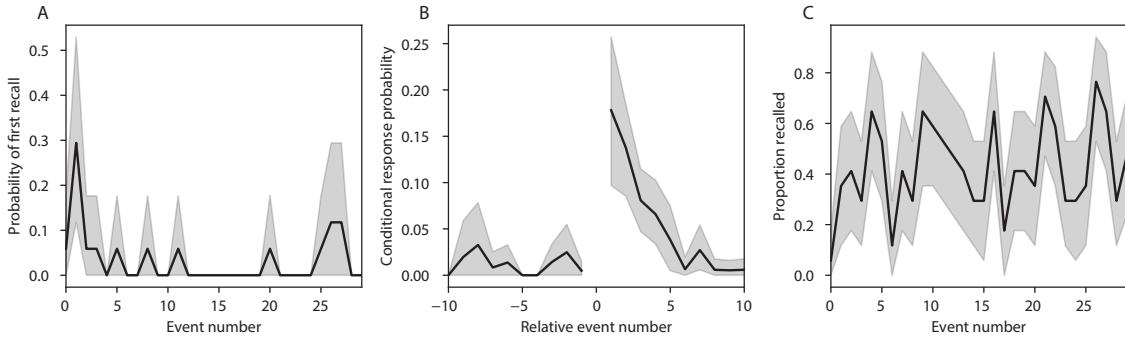


Figure 3: Naturalistic extensions of classic list-learning memory analyses. A. The probability of first recall as a function of the serial position of the event in the video. B. The probability of recalling each event, conditioned on having most recently recalled the event *lag* events away in the video.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

236 SEM: 0.029), and that participants with higher temporal clustering scores tended to perform better
 237 according to both Chen et al. (2017)'s hand-annotated memory scores (Pearson's $r(15) = 0.62$, $p =$
 238 0.008) and our model's estimate (Pearson's $r(15) = 0.54$, $p = 0.024$). Semantic clustering measures
 239 the extent to which participants cluster their recall responses according to semantic similarity.
 240 We found that participants tended to recall semantically similar video events together (mean
 241 clustering score: 0.787, SEM: 0.018), and that semantic clustering score was also related to both
 242 hand-annotated (Pearson's $r(15) = 0.65$, $p = 0.004$) and model-derived (Pearson's $r(15) = 0.63$, $p =$
 243 0.007) memory performance.

244 Statistical models of memory studies often treat recall success as binary (in other words, an
 245 item either was or was not recalled), or occasionally categorical (e.g., to distinguish familiarity
 246 from recollection; Yonelinas et al., 2002). Such approaches are tenable in classical list-learning or
 247 recognition memory paradigms, as the presented stimuli tend to be very simple (e.g., a sequence
 248 of individual words or items). However, memory for naturalistic experiences is much more
 249 nuanced. For example, certain aspects of an experience might be correctly remembered at varying
 250 levels of detail, or distorted, or forgotten entirely. Further, each remembering is itself a richly
 251 structured phenomenon. Our framework produces a content-based model of individual video
 252 and recall events by projecting the dynamic content of the video and participants' recalls into a

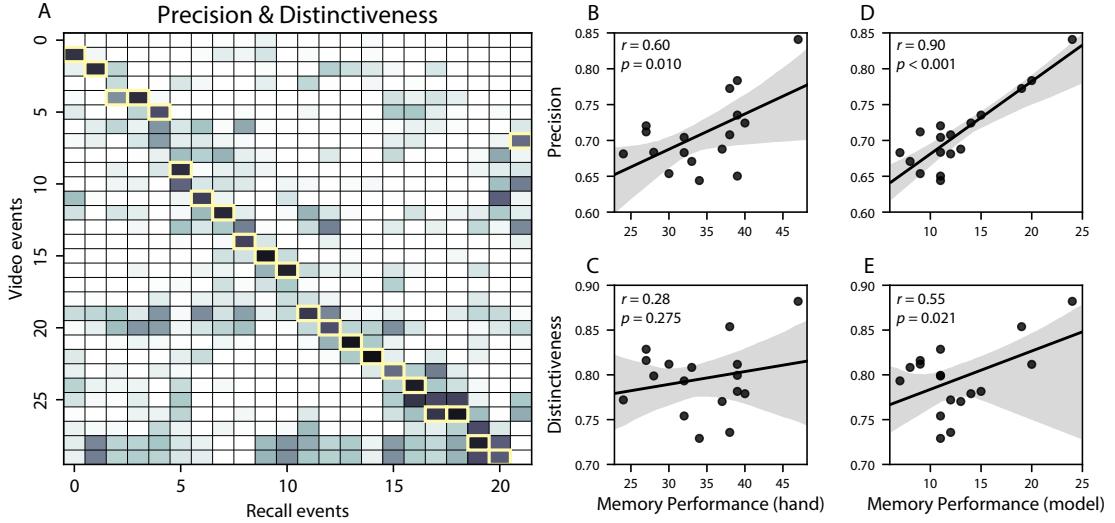


Figure 4: Novel content-based metrics of naturalistic memory: precision and distinctiveness. **A.** The video-recall correlation matrix for a representative participant (17). The yellow boxes highlight the maximum correlation in each column. The example participant's overall precision score was computed as the average across correlation values in the yellow boxes. Their distinctiveness score was computed as the average (over recall events) of 1 minus the average correlation between each recall event and all other recall events that do not display a box in the same row. **B.** The (Pearson's) correlation between precision and hand-annotated memory performance. **C.** The correlation between distinctiveness and hand-annotated memory performance. **D.** The correlation between precision and the number of video events successfully recalled, as determined by our model. **E.** The correlation between distinctiveness and the number of video events successfully recalled, as determined by our model.

shared topic space. This allows for direct, quantitative comparisons 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, *continuous* metrics for analyzing naturalistic memory: *precision* and *distinctiveness*. Precision is intended to capture the “completeness” of recall, or how fully the presented content was recapitulated in memory. We define a recall event’s precision as the maximum correlation between the topic proportions of that recall event and any video event (Fig. 4). A second novel metric we introduce here is *distinctiveness*, which is intended to capture the “specificity” of recall. In other words, distinctiveness quantifies the extent to which a given recalled event reflects the most similar presented event moreso than it does other presented events. To compute a recall event’s distinctiveness, we first identify the

263 video event to which its topic vector is most strongly correlated. We then define distinctiveness
264 as one minus the average correlation between the given recall event and all *other* video events.
265 In addition to individual events, one may also use these metrics to describe each participant's
266 overall performance by averaging across a participant's event-wise precision or distinctiveness
267 scores. Participants whose recall events are more veridical descriptions of what happened in
268 the video event will presumably have higher precision scores. We find that, across participants,
269 higher precision scores are positively correlated with both hand-annotated memory performance
270 (as collected by Chen et al., 2017; Pearson's $r(15) = 0.60, p = 0.010$) and the number of video events
271 successfully remembered, as determined by our model (Pearson's $r(15) = 0.90, p < 0.001$). We also
272 hypothesized that participants who recounted events in a more distinctive way would display
273 better overall memory. We find that participants' distinctiveness scores were positively correlated
274 with our model's estimated number of recall events (Pearson's $r(15) = 0.55, p = 0.021$). However,
275 we found no evidence that distinctiveness scores were correlated with hand-annotated memory
276 performance (Pearson's $r(15) = 0.28, p = 0.275$). We elaborate on this potential discrepancy in the
277 *Discussion* section.

278 Further intuition for the behaviors captured by these two metrics may be gained by directly
279 examining the content of the video and recalls our framework models. In Figure 5, we contrast
280 recalls for the same video event (event 22) from two participants: one with a high precision score
281 (P17), the other with a low precision score (P6). From the HMM-identified event boundaries,
282 we recovered the set of annotations describing the content of an example video event (Fig. 5B),
283 and divided them into different color-coded sections for each action or feature described. We
284 then similarly recovered the set of sentences comprising the corresponding recall event for each
285 of the two example participants. Because the recall sliding windows overlap heavily, and each
286 recall event spans multiple recall timepoints (i.e., windows), we have stripped any sentences from
287 the beginning and end that describe earlier or later video events for the sake of readability. In
288 other words, Fig. 5C shows a subset of the full recall event text, comprising sentences between
289 the first and last descriptions of content from the example video event. We then colored all words
290 describing actions and features coded in panel B by their corresponding color. Visual comparison

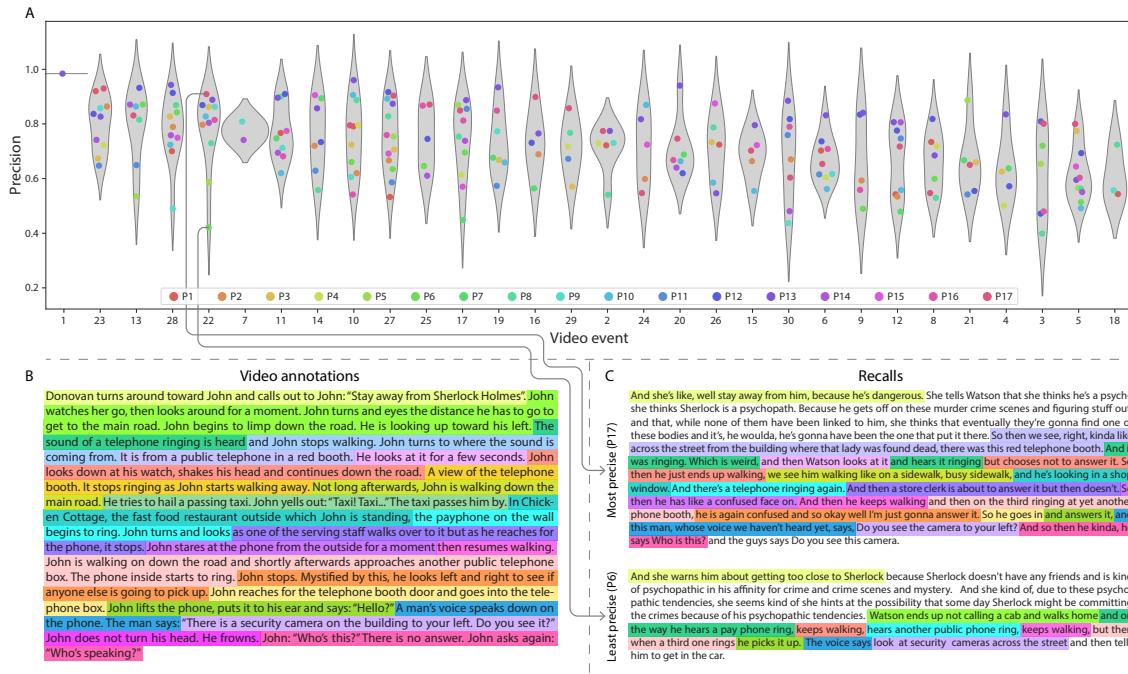


Figure 5: Precision metric reflects completeness of recall. **A.** Recall precision by video event. Grey violin plots display kernel density estimates for the distribution of recall precision scores for a single video event. Colored dots within each violin plot represent individual participants' recall precision for the given event. Video events are ordered along the *x*-axis by the average precision with which they were remembered. **B.** The set of "Narrative Details" video annotations (generated by Chen et al., 2017) for scenes comprising an example video event (22) identified by the HMM. Each action or feature is highlighted in a different color. **C.** A subset of the sentences comprising the most precise (P17) and least precise (P6) participants' recalls of video event 22. Descriptions of specific actions or features reflecting those highlighted in panel B are highlighted in the corresponding color.

291 of these example transcripts reveals that the more precise recall captures more of the video event's
 292 content, and with more detail.

293 Figure 6 similarly contrasts two example participants' recalls for a common video event (event
 294 19) to illustrate the tangible differences between high and low distinctiveness scores. Here, we
 295 have extracted the full set of sentences comprising the most distinctive recall event (P13) and least
 296 distinctive recall event (P11) matched to the example video event (Fig. 6C). We also extracted the
 297 annotations for the example video event, as well as those from each other video event whose content
 298 the example participants' single recall events described (Fig. 6B). We then shaded the annotation
 299 text for each video event with a different color, and shaded each word of the example participants'

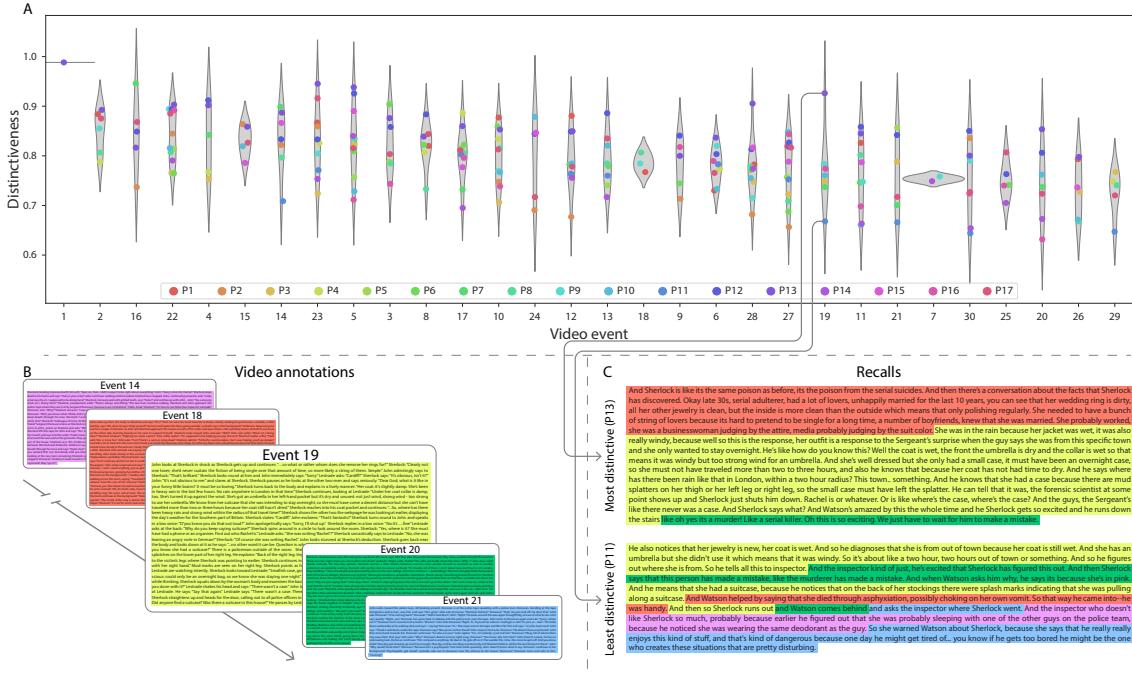


Figure 6: Distinctiveness metric reflects specificity of recall. A. Recall distinctiveness by video event. Kernel density estimates for each video event’s distribution of recall distinctiveness scores, analogous to Fig. 5A. B. The sets of “Narrative Details” video annotations (generated by Chen et al., 2017) for scenes comprising video events described by the example participants in panel C. Each event’s text is highlighted in a different color. C. The sentences comprising the most distinctive (P13) and least distinctive (P11) participants’ recalls of video event 19. Sections of recall describing each video event in panel B are highlighted with the corresponding color.

recall text by the color of the video event it describes. The majority of the most distinctive recall
300 event text describes video event 19's content, with the first five and last one sentence describing
301 the video events immediately preceding and succeeding the current one, respectively. In contrast,
302 the least distinctive recall of video event 19 blends the content from five separate video events,
303 does not transition between them in order, and often combines descriptions of two video events'
304 content in the same sentence.

306 The prior analyses leverage the correspondence between the 100-dimensional topic proportion
307 matrices for the video and participants' recalls to characterize recall. However, it is difficult to
308 gain deep insights into the content of (or relationships between) experiences and memories solely
309 by examining these topic proportions (e.g., Figs. 2A, D) or the corresponding correlation matrices

(Figs. 2B, E, S4). And while we can directly examine the original text underlying these topic vectors (e.g., Figs. 5, 6) to show how relationships between them reflect real-world behavior, this comparison becomes prohibitively cumbersome at larger timescales. To visualize the time-varying high-dimensional content in a more intuitive way (Heusser et al., 2018b), we projected the topic proportions matrices onto a two-dimensional space using Uniform Manifold Approximation and Projection (UMAP; McInnes et al., 2018). In this lower-dimensional space, each point represents a single video or recall event, and the distances between the points reflect the distances between the events' associated topic vectors (Fig. 7). In other words, events that are nearer to each other in this space are more semantically similar, and those that are farther apart are less so.

Visual inspection of the video and recall topic trajectories reveals a striking pattern. First, the topic trajectory of the video (which reflects its dynamic content; Fig. 7A) is captured nearly perfectly by the averaged topic trajectories of participants' recalls (Fig. 7B). To assess the consistency of these recall trajectories across participants, we asked: given that a participant's recall trajectory had entered a particular location in the reduced topic space, could the position of their *next* recalled event be predicted reliably? For each location in the reduced topic space, we computed the set of line segments connecting successively recalled events (across all participants) that intersected that location (see *Methods* for additional details). We then computed (for each location) the distribution of angles formed by the lines defined by those line segments and a fixed reference line (the *x*-axis). Rayleigh tests revealed the set of locations in topic space at which these across-participant distributions exhibited reliable peaks (blue arrows in Fig. 7B reflect significant peaks at $p < 0.05$, corrected). We observed that the locations traversed by nearly the entire video trajectory exhibited such peaks. In other words, participants exhibited similar trajectories that also matched the trajectory of the original video (Fig. 7C). This is especially notable when considering the fact that the number of events participants recalled (dots in Fig. 7C) varied considerably across people, and that every participant used different words to describe what they had remembered happening in the video. Differences in the numbers of remembered events appear in participants' trajectories as differences in the sampling resolution along the trajectory. We note that this framework also provides a means of disentangling classic "proportion recalled" measures (i.e., the proportion

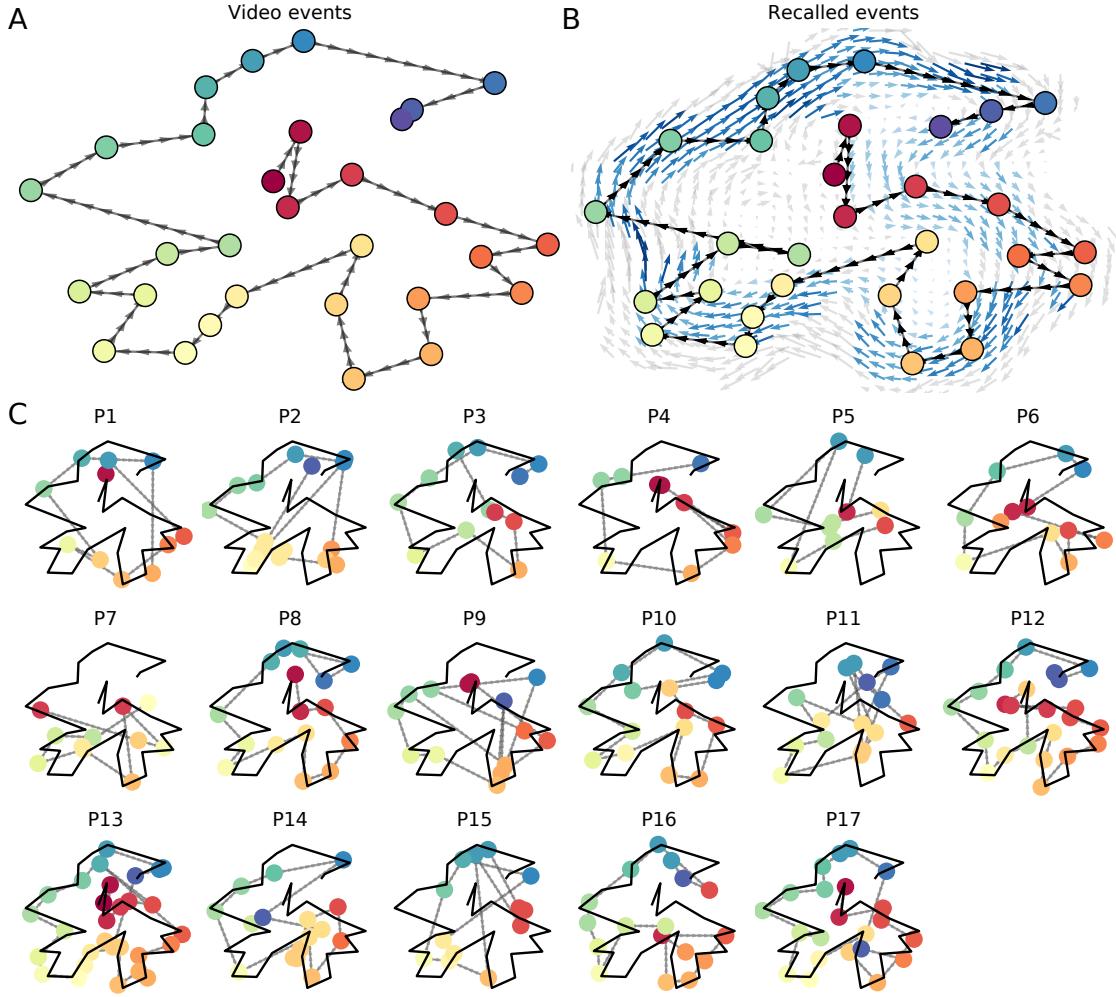


Figure 7: Trajectories through topic space capture the dynamic content of the video and recalls. All panels: the topic proportion matrices have been projected onto a shared two-dimensional space using UMAP. **A.** The two-dimensional topic trajectory taken by the episode of *Sherlock*. Each dot indicates 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. Here, events (dots) are colored by their matched video event (Panel A).

338 of video events described in participants' recalls) from participants' abilities to recapitulate the
339 overall unfolding of the original video's content (i.e., the similarity between the shapes of the
340 original video trajectory and that defined by each participant's recounting of the video).

341 In addition to the more "holistic" measure of memory described in the previous section, our
342 framework also affords the ability to drill down to individual words and quantify how each word
343 relates to the memorability of each event. The results displayed in Figures 3C and 5A suggest that
344 certain events were remembered better than others. Given this, we next asked whether the
345 events were generally remembered well or poorly tended to reflect particular content. Because
346 our analysis framework projects the dynamic video content and participants' recalls into a shared
347 space, and because the dimensions of that space represent topics (which are, in turn, sets of
348 weights over known words in the vocabulary), we are able to recover the weighted combination
349 of words that make up any point (i.e., topic vector) in this space. We first computed the average
350 precision with which participants recalled each of the 30 video events (Fig. 8A; note that this result
351 is analogous to a serial position curve created from our continuous recall quality metric). We
352 then computed a weighted average of the topic vectors for each video event, where the weights
353 reflected how reliably each event was recalled. To visualize the result, we created a "wordle"
354 image (Mueller et al., 2018) where words weighted more heavily by better-remembered topics
355 appear in a larger font (Fig. 8B, green box). Across the full video, content that reflected topics
356 necessary to convey the central focus of the video (e.g., the names of the two main characters,
357 "Sherlock" and "John," and the address of a major recurring location, "221B Baker Street") were
358 best remembered. An analogous analysis revealed which themes were poorly remembered. Here
359 in computing the weighted average over events' topic vectors, we weighted each event in *inverse*
360 proportion to how well it was remembered (Fig. 8B, red box). The least well-remembered video
361 content reflected information not necessary to later convey a general summary of the video, such
362 as the proper names of relatively minor characters (e.g., "Mike," "Molly," and "Lestrade") and
363 locations (e.g., "St. Bartholomew's Hospital").

364 A similar result emerged from assessing the topic vectors for individual video and recall events
365 (Fig. 8C). Here, for each of the three best- and worst-remembered video events, we have constructed

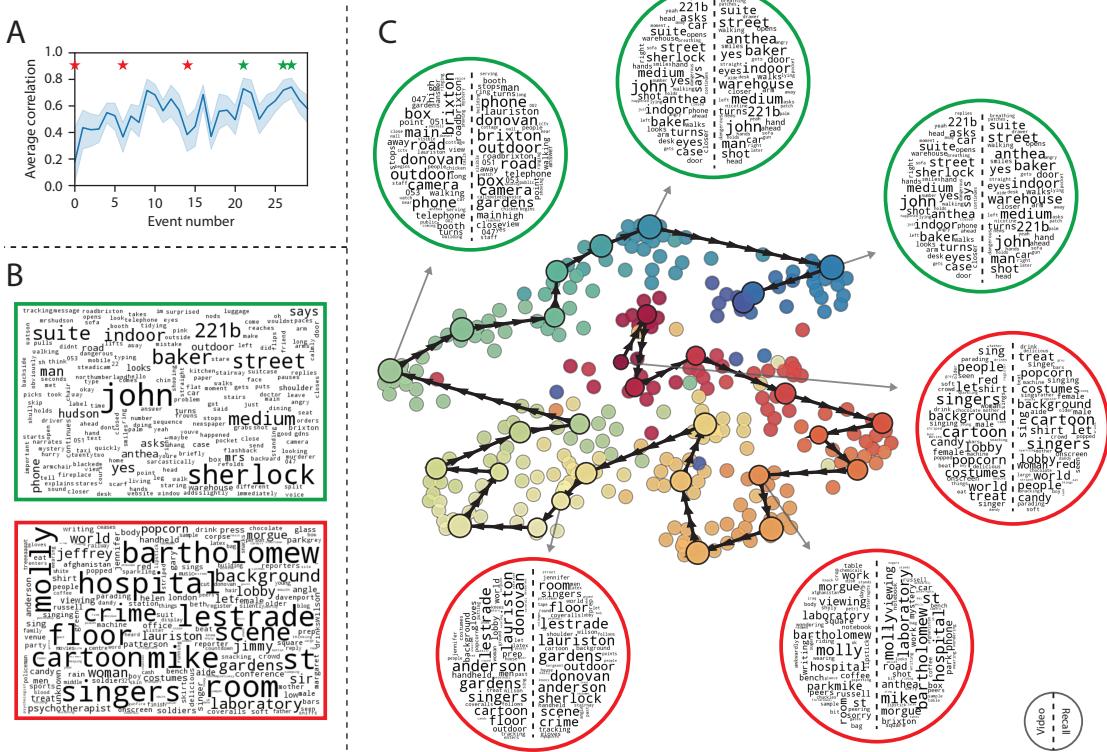


Figure 8: Language used in the most and least memorable events. **A.** Average precision (video event-recall event topic vector correlation) across participants for each video event. Error bars denote bootstrap-derived across-participant 95% confidence intervals. The stars denote the three best-remembered 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 7. 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 7A. 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.

366 two wordles: one from the original video event's topic vector (left) and a second from the average
367 recall topic vector for that event (right). The three best-remembered events (circled in green)
368 correspond to scenes integral to the central plot-line: a mysterious figure spying on John in a
369 phone booth; John meeting Sherlock at Baker St. to discuss the murders; and Sherlock laying
370 a trap to catch the killer. Meanwhile, the three worst-remembered events (circled in red) reflect
371 scenes that are non-essential to summarizing the narrative's structure: the video of singing cartoon
372 characters participants viewed in an introductory clip prior to the main episode; John asking Molly
373 about Sherlock's habit of over-analyzing people; and Sherlock noticing evidence of Anderson's
374 and Donovan's affair.

375 The results thus far inform us about which aspects of the dynamic content in the episode partic-
376 ipants watched were preserved or altered in participants' memories. We next carried out a series
377 of analyses aimed at understanding which brain structures might facilitate these preservations
378 and transformations between the external world and memory. In the first analysis, we sought
379 to identify brain structures that were sensitive to the dynamic unfolding of the video's content,
380 as characterized by its topic trajectory. We used a searchlight procedure to identify clusters of
381 voxels whose activity patterns displayed a proximal temporal correlation structure (as participants
382 watched the video) matching that of the original video's topic proportions (Fig. 9A; see *Methods* for
383 additional details). In a second analysis, we sought to identify brain structures whose responses
384 (during video viewing) reflected how each participant would later structure their recounting of the
385 video. We used an analogous searchlight procedure to identify clusters of voxels whose proximal
386 temporal correlation matrices matched that of the topic proportions for each individual's recall
387 (Figs. 9B; see *Methods* for additional details). To ensure our searchlight procedure identified re-
388 gions *specifically* sensitive to the temporal structure of the video or recalls (i.e., rather than those
389 with a temporal autocorrelation length similar to that of the video/recalls), we performed a phase
390 shift-based permutation correction (see *Methods* for additional details). As shown in Figure 9C, the
391 video-driven searchlight analysis revealed a distributed network of regions that may play a role in
392 processing information relevant to the narrative structure of the video. Similarly, the recall-driven
393 searchlight analysis revealed a second network of regions (Fig. 9D) that may facilitate a person-

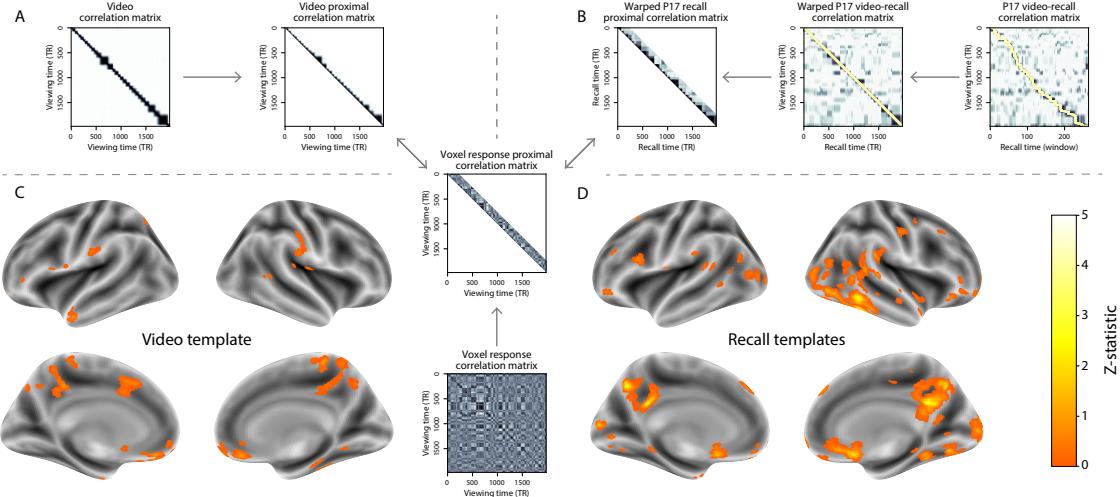


Figure 9: Brain structures that underlie the transformation of experience into memory. **A.** We isolated the proximal diagonals from the upper triangle of the video correlation matrix, and applied this same diagonal mask to the voxel response correlation matrix for each cube of voxels in the brain. We then searched for brain regions whose activation timeseries consistently exhibited a similar proximal correlational structure to the video model, across participants. **B.** We used dynamic time warping (Berndt and Clifford, 1994) to align each participant's recall timeseries to the TR timeseries of the video. We then applied the same diagonal mask used in Panel A to isolate the proximal temporal correlations and searched for brain regions whose activation timeseries for an individual consistently exhibited a similar proximal correlational structure to each individual's recall. **C.** We identified a network of regions sensitive to the narrative structure of participants' ongoing experience. The map shown is thresholded at $p < 0.05$, corrected. **D.** We also identified a network or regions sensitive to how individuals would later structure the video's content in their recalls. The map shown is thresholded at $p < 0.05$, corrected.

394 specific transformation of one's experience into memory. In identifying regions whose responses
 395 to ongoing experiences reflect how those experiences will be remembered later, this latter analysis
 396 extends classic *subsequent memory analyses* (e.g., Paller and Wagner, 2002) to domain of naturalistic
 397 stimuli.

398 The searchlight analyses described above yielded two distributed networks of brain regions,
 399 whose activity timecourses mirrored to the temporal structure of the video (Fig. 9C) or participants'
 400 eventual recalls (Fig. 9D). We next sought to gain greater insight into the structures and functional
 401 networks our results reflected. To accomplish this, we performed an additional, exploratory
 402 analysis using Neurosynth (Yarkoni et al., 2011). Given an arbitrary statistical map as input,
 403 Neurosynth performs a massive automated meta-analysis, returning a ranked list of terms reported

⁴⁰⁴ in papers with similar significance maps. We ran Neurosynth on the significance maps for the video-
⁴⁰⁵ and recall-driven searchlight analyses. These maps, along with the 10 terms with maximally similar
⁴⁰⁶ meta-analysis images identified by Neurosynth are shown in Figure 10.

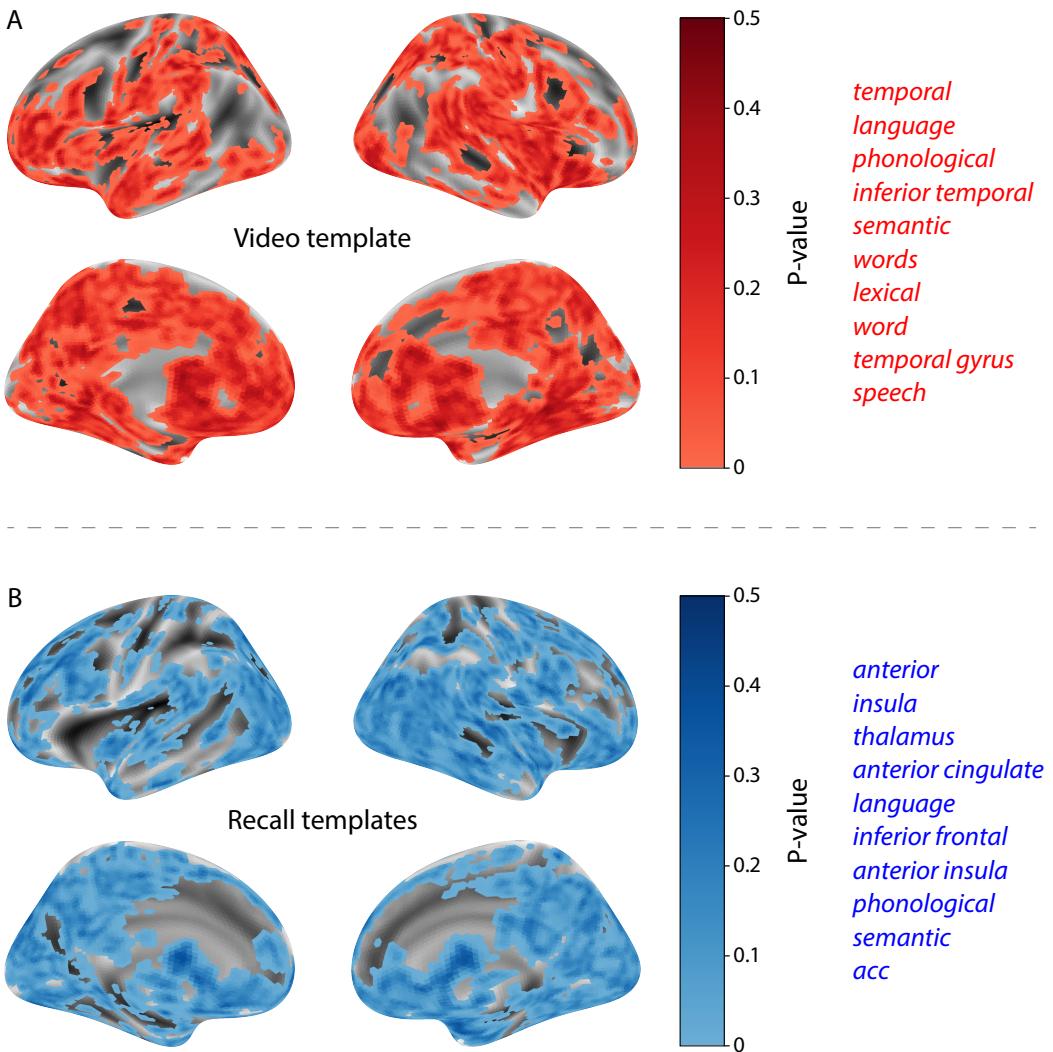


Figure 10: Decoding distributed statistical maps via Neurosynth meta-analyses. **A.** Video-searchlight significance and top 10 decoded terms. We constructed a map of the permutation-derived p -values for the video-driven searchlight analysis (Fig. 9A, C) at each voxel with a positive permutation-derived z -score. The top 10 terms decoded from this significance map are shown in red. **B.** Recall-searchlight significance and top 10 decoded terms. We constructed a map of the permutation-derived p -values for the recall-driven searchlight analysis (Fig. 9A, C) at each voxel with a positive permutation-derived z -score. The top 10 terms decoded from this significance map are shown in blue.

407 **Discussion**

408 Our work casts remembering as reproducing (behaviorally and neurally) the topic trajectory, or
409 shape, of an experience. This view draws inspiration from prior work aimed at elucidating
410 the neural and behavioral underpinnings of how we process dynamic naturalistic experiences
411 and remember them later. One approach to identifying neural responses to naturalistic stimuli
412 (including experiences) entails building a model of the stimulus and searching for brain regions
413 whose responses are consistent with the model. In prior work, a series of studies from Uri Hasson's
414 group (Lerner et al., 2011; Simony et al., 2016; Chen et al., 2017; Baldassano et al., 2017; Zadbood
415 et al., 2017) have extended this approach with a clever twist: rather than building an explicit
416 stimulus model, these studies instead search for brain responses (while experiencing the stimulus)
417 that are reliably similar across individuals. So called *inter-subject correlation* (ISC) and *inter-subject*
418 *functional connectivity* (ISFC) analyses effectively treat other people's brain responses to the stimulus
419 as a "model" of how its features change over time. By contrast, in our present work, we use topic
420 models to construct an explicit content model directly from the stimulus (i.e., the topic trajectory
421 of the video). Projecting each participant's recall into a space shared by both the stimulus and
422 other participants then allows us to compare recalls both directly to the stimulus and to each other.
423 Similarly, prior work introducing the use of HMMs to discover latent event structure in naturalistic
424 stimuli and recall (Baldassano et al., 2017) used between-subjects cross-validation to identify event
425 boundaries shared across participants, and between stimulus and recall. Our framework allows
426 us to break from the restriction of a common, shared event-timeseries and identify the unique
427 *resolution* of each participant's recall event structure, and how that may differ from the video and
428 that of other participants.

429 Word embedding models are a rapidly growing area of machine learning research. Early ap-
430 proaches including latent semantic analysis (Landauer and Dumais, 1997) use word co-occurrence
431 statistics (i.e., how often pairs of words occur in the same documents contained in the corpus) to
432 derive a unique feature vector for each word. The feature vectors are constructed so that words
433 that co-occur more frequently have feature vectors that are closer (in Euclidean distance). Related

434 approaches, such as latent dirichlet allocation (Blei et al., 2003) attempt to explicitly model the
435 underlying causes of word co-occurrences by automatically identifying the set of themes or topics
436 reflected across the documents in the corpus. More recent work on these types of semantic mod-
437 els, including word2vec (Mikolov et al., 2013), the Universal Sentence Encoder (Cer et al., 2018),
438 GPT-2 (Radford et al., 2019), and GTP-3 (Brown et al., 2020) use deep neural networks to attempt
439 to identify the deeper conceptual representations underlying each word. Despite the growing
440 popularity of more sophisticated deep learning-based embedding models, here we leverage latent
441 dirichlet allocation (i.e., topic modelling) to embed video and recall text. This decision was mo-
442 tivated by several factors. First, topic models capture the *essence* of a text passage devoid of the
443 specific set and order of words used. This was an important feature of our model since different
444 people may accurately recall a scene using very different language. Second, words can mean dif-
445 ferent things in different contexts (e.g. “bat” may be the act of hitting a baseball, the object used for
446 that action, or as a flying mammal). Topic models are robust to this, allowing words to exist as part
447 of multiple topics. Last, topic models provide a straightforward means of recovering the weights
448 for the particular words comprising a topic, enabling straightforward interpretation of an event’s
449 contents (e.g. Fig. 8). Other models such as the Universal Sentence Encoder, GPT-2, and GPT-3
450 offer context-sensitive encoding of text passages, but the encoding space is complex and non-linear,
451 and thus recovering the original words used to fit the model is not straightforward. However, it is
452 worth pointing out that our general framework is divorced from the particular choice of language
453 model. Moreover, many of the aspects of our framework could be swapped out for other choices.
454 For example, the language model, the timeseries segmentation model and the video-recall match-
455 ing function could all be customized to suit a particular question space or application. Indeed for
456 some questions, recovery of the particular words used to describe a memory may not be necessary,
457 and thus other text-modeling approaches (including the deep learning-based embedding models
458 described above) may be preferable. Future work will explore the influence of particular model
459 choices on the framework’s efficacy.

460 In extending classical free recall analyses to our naturalistic memory framework, we recovered
461 two patterns of recall dynamics central to list-learning studies: a heightened probability of initiating

recall with the first presented “item” (in our case, video events; Fig. 3A) and a strong bias toward transitioning from recalling a given event to recalling the one immediately following it (Fig. 3B). However, equally noteworthy are the typical free recall results *not* recovered in these analyses, as each highlights a fundamental difference between the list-learning paradigm and naturalistic memory paradigms like the one employed in the present study. The most noticeable departure from hallmark free recall dynamics in these findings is the apparent lack of a serial position effect in Figure 3C, which instead shows greater and lesser recall probabilities for events distributed across the video. Stimuli in free recall experiments most often comprise lists of simple, common words, presented to participants in a random order. (In fact, numerous word pools have been developed based on these criteria; e.g., Friendly et al., 1982). These stimulus qualities enable two assumptions that are central to word list analyses, but frequently do not hold for real-world experiences. First, researchers conducting list-learning studies may assume that the content at each presentation index is essentially equal, and does not possess attributes that would render it, on average, more or less memorable than others. Such is rarely the case with real-world experiences or experiments meant to approximate them, and the effects of both intrinsic and observer-dependent factors on stimulus memorability are well established (for review see Chun and Turk-Browne, 2007; Bylinskii et al., 2015; Tyng et al., 2017). Second, the random ordering of list items ensures that (across participants, on average) there is no relationship between the thematic similarity of individual stimuli and their presentation positions—in other words, two successively presented items are no more likely to be highly semantically similar than they are to be highly dissimilar. In most cases, the exact opposite is true of real-world episodes. Our internal thoughts, our actions, and the physical state of the world around us all tend to follow a direct (often causal) progression. As a result, each moment of our experience tends to be inherently more similar to surrounding moments than to those in the distant past or future. Memory literature has termed this strong temporal autocorrelation “context,” and in various media that depict real-world events (e.g., movies or written stories), we recognize it as a *narrative structure*. While a random word list (by definition) has no such structure, the logical progression between ideas and actions in a naturalistic stimulus prompts the rememberer to recount presented events in order, starting with the beginning. This tendency is

490 reflected in our findings' second departure from typical free recall dynamics: a lack of increased
491 probability of first recall for end-of-sequence events (Fig. 3A).

492 Because they disregard presentation order-dependent variability in the stimulus content, analyses
493 such as those in Figure 3 enable a more sensitive analysis of presentation order-dependent
494 temporal dynamics in free recall. Yet by the same token, they paint a wholly incomplete picture of
495 memory for naturalistic episodes. In an attempt to address this shortcoming, we have developed a
496 framework in the present study that characterizes the explicit semantic content of the stimulus and
497 subsequent recalls. However, sensitivity to stimulus and recall content introduces a new challenge:
498 distinguishing between levels of recall quality for a stimulus (e.g., an event) that is considered to
499 have been "remembered." When modeling memory in an experimental setting, recall quality for
500 individual events is often cast as binary (e.g., a given list item was simply either remembered or
501 not remembered). Various models of memory (e.g., Yonelinas, 2002) attempt to improve upon this
502 by including confidence ratings, rendering this binary judgement instead categorical. To better
503 evaluate naturalistic memory quality, we introduce a continuous metric (*precision*), which reflects
504 the level of completeness of a participant's recall for a feature-rich experience. Additionally, recall
505 quality for a single event is typically assessed independently from that for all other events (e.g., it
506 is difficult to "compare" a participant's binary recall success for list item 1 to that of list item 10).
507 The second novel metric we introduce (*distinctiveness*) is based on analyzing of the correlational
508 structure of an individual's full set of recall events, and reflects the specificity of their memory for
509 a single experienced event. We find that both of these metrics relate to the overall number of video
510 events participants successfully recalled, and that our precision metric additionally relates to Chen
511 et al. (2017)'s hand-annotated memory memory scores.

512 We did not find evidence that participants' average recall distinctiveness was related to their
513 hand-annotated memory scores computed by Chen et al. (2017). One possible explanation is that,
514 in hand-scoring each participant's verbal recall for each of 50 (manually-delimited) scenes, "[a]
515 scene was counted as recalled if the participant described any part of the scene" (Chen et al.,
516 2017). In other words, both an extensive description of a scene's content and a brief mention of
517 some subset of its content were (binarily) considered equally successful recalls. By contrast, we

518 identify the event structure in participants' recalls in an unsupervised manner, independent of
519 the video event-timeseries, prior to mapping between video and recall content. Our HMM-based
520 event-segmentation produces boundaries between timepoints where the topic proportions shift in
521 a substantial way, and because a small handful of words is unlikely to contribute significantly to
522 the topic proportions for any sliding window, such brief scene descriptions will most often not
523 result in a sufficiently large shift in the resulting topic proportions for the HMM to identify an
524 event boundary. Instead, they will be grouped with a neighboring event, consequently lowering
525 that event's distinctiveness score and by extension, the participant's overall distinctiveness score.
526 This is in essence the qualitative difference between distinctive and indistinctive recall, and reflects
527 the comparison shown in Figure 6C. Intriguingly, prior studies show that pattern separation, or the
528 ability to cleanly discriminate between similar experiences, is impaired in many cognitive disorders
529 as well as natural aging (Stark et al., 2010; Yassa et al., 2011; Yassa and Stark, 2011). Future work
530 might explore whether and how these metrics compare between cognitively impoverished groups
531 and healthy controls.

532 In the analyses outlined in Figure 9, we identified two networks of brain regions whose re-
533 sponses during video viewing were consistent with the temporal structure of the video and recall
534 topic trajectories, respectively. The network identified by the video trajectory analysis included the
535 ventromedial prefrontal cortex, left anterior temporal lobe, superior parietal and dorsal anterior
536 cingulate cortex. The network from the video-recall trajectory analysis also included the ventro-
537 medial prefrontal and superior parietal cortices, in addition to the posterior medial cortex (PMC)
538 and the inferior temporal regions. Notably, Chen et al. (2017) also observed the PMC in a number
539 of analyses including one that searched for regions whose activity patterns during encoding were
540 reinstated during free recall. The PMC has been consistently identified in studies involving stimuli
541 with meaningfully structured events (Cohn-Sheely and Ranganath, 2017). Further, the PMC is
542 part of the "posterior medial" system, a network of brain regions thought to represent situation
543 models (Zacks et al., 2007) in support of memory, spatial navigation and social cognition (Ran-
544 ganath and Ritchey, 2012). Given that we constructed our video-recall searchlight model to capture
545 temporal structure in the episode's semantic content (and how one's later recall aligns with that

546 structure), we speculate that the PMC may play a role in constructing mnemonic events from
547 meaningfully structured experiences.

548 Decoding the associated significance maps with Neurosynth revealed two intriguing results.
549 First, the top 10 terms returned for the video-driven searchlight significance map were centered
550 around themes of language and semantic meaning (Fig. 10A). In other words, the voxels identified
551 as more reflective of the video content's temporal structure (i.e., voxels with lower permutation
552 correction-derived p -values), as defined by our model, were most likely to be reported as active in
553 studies focused on the the neural underpinnings of semantic processing. This finding is interesting,
554 as our model specifically captures the temporal structure of the video's *semantic* content (e.g., as
555 opposed to that of the visual, auditory, or affective content). This suggests that the network of
556 structures displayed in Figure 9C may play a roll in processing the evolving semantic content of
557 ongoing experiences.

558 Our second searchlight analysis identified a partially overlapping network of regions (Fig. 9D)
559 whose patterns of activity as participants viewed the video reflected the idiosyncratic structure of
560 each individual's later recalls. The associated significance map yielded a set of Neurosynth terms
561 that primarily reflected names of specific structural regions (such as "thalamus," "anterior insula,"
562 "anterior cingulate" and "inferior frontal"; Fig. 10B). Interestingly, these regions share membership
563 in a common, large-scale functional network (termed the "salience network") involved in detect-
564 ing and processing affective cues. In particular, the latter three regions have been implicated in
565 functions relevant to assigning personal meaning to an experience, including: ascribing subjective
566 value to raw, sensory input (Medford and Critchley, 2010); modulating semantic and phonological
567 processing in response to personally salient stimuli (Kelly et al., 2007); and directing and reallo-
568 cating attention and working memory resources towards the most relevant stimuli (Menon and
569 Uddin, 2010). This suggests that the network of structures displayed in Figure 9D may be play a roll
570 in transforming and restructuring ongoing experiences through the lens of one's prior experience
571 and subjective emotions as they are encoded in memory.

572 Our work has broad implications for how we characterize and assess memory in real-world
573 settings, such as the classroom or physician's office. For example, the most commonly used

574 classroom evaluation tools involve simply computing the proportion of correctly answered exam
575 questions. Our work indicates that this approach is only loosely related to what educators might
576 really want to measure: how well did the students understand the key ideas presented in the
577 course? Under this typical framework of assessment, the same exam score of 50% could be
578 ascribed to two very different students: one who attended the full course but struggled to learn
579 more than a broad overview of the material, and one who attended only half of the course but
580 understood the material perfectly. Instead, one could apply our computational framework to build
581 explicit content models of the course material and exam questions. This approach would provide
582 a more nuanced and specific view into which aspects of the material students had learned well
583 (or poorly). In clinical settings, memory measures that incorporate such explicit content models
584 might also provide more direct evaluations of patients' memories.

585 Methods

586 Experimental design and data collection

587 Data were collected by Chen et al. (2017). In brief, participants ($n = 22$) viewed the first 48 minutes
588 of "A Study in Pink", the first episode of the BBC television series *Sherlock*, while fMRI volumes
589 were collected (TR = 1500 ms). Participants were pre-screened to ensure they had never seen any
590 episode of the show before. The stimulus was divided into a 23 min (946 TR) and a 25 min (1030 TR)
591 segment to mitigate technical issues related to the scanner. After finishing the clip, participants
592 were instructed to (quoting from Chen et al., 2017) "describe what they recalled of the [episode]
593 in as much detail as they could, to try to recount events in the original order they were viewed
594 in, and to speak for at least 10 minutes if possible but that longer was better. They were told that
595 completeness and detail were more important than temporal order, and that if at any point they
596 realized they had missed something, to return to it. Participants were then allowed to speak for
597 as long as they wished, and verbally indicated when they were finished (e.g., 'I'm done')." Five
598 participants were dropped from the original dataset due to excessive head motion (2 participants),

599 insufficient recall length (2 participants), or falling asleep during stimulus viewing (1 participant),
600 resulting in a final sample size of $n = 17$. For additional details about the experimental procedure
601 and scanning parameters, see Chen et al. (2017). The experimental protocol was approved by
602 Princeton University's Institutional Review Board.

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

609 The video stimulus was divided into 1,000 fine-grained “scenes” and annotated by an inde-
610 pendent coder. For each of these 1,000 scenes, the following information was recorded: a brief
611 narrative description of what was happening, the location where the scene took place, whether
612 that location was indoors or outdoors, the names of all characters on-screen, the name(s) of the
613 character(s) in focus in the shot, the name(s) of the character(s) currently speaking, the camera
614 angle of the shot, a transcription of any text appearing on-screen, and whether or not there was
615 music present in the background. Each scene was also tagged with its onset and offset time, in
616 both seconds and TRs.

617 **Data and code availability**

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

620 **Statistics**

621 All statistical tests performed in the behavioral analyses were two-sided. All statistical tests per-
622 formed in the neural data analyses were two-sided, except for the permutation-based thresholding,
623 which was one-sided. In this case, we were specifically interested in identifying voxels whose ac-

624 tivation time series reflected the temporal structure of the video and recall trajectories to a *greater*
625 extent than that of the phase-shifted trajectories.

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

627 **Topic modeling**

628 The input to the topic model we trained to characterize the dynamic content of the video comprised
629 998 hand-generated annotations of short (mean: 2.96s) scenes spanning the video clip (Chen
630 et al., 2017 generated 1000 annotations total; we removed two annotations referring to a break
631 between the first and second scan sessions, during which no fMRI data was collected). We
632 concatenated the text for all of the annotated features within each segment, creating a “bag of
633 words” describing each scene and performed some minor preprocessing (e.g., stemming possessive
634 nouns and removing punctuation). We then re-organized the text descriptions into overlapping
635 sliding windows spanning (up to) 50 scenes each. In other words, we estimated the “context”
636 for each scene using the text descriptions of the preceding 25 scenes, the present scene, and the
637 following 24 scenes. To model the context for scenes near the beginning of the video (i.e., within
638 25 scenes of the beginning or end), we created overlapping sliding windows that grew in size
639 from one scene to the full length. We also tapered the sliding window lengths at the end of the
640 video, whereby scenes within fewer than 24 scenes of the end of the video were assigned sliding
641 windows that extended to the end of the video. This procedure ensured that each scene’s content
642 was represented in the text corpus an equal number of times.

643 We trained our model using these overlapping text samples with `scikit-learn` (version 0.19.1;
644 Pedregosa et al., 2011), called from our high-dimensional visualization and text analysis software,
645 `HyperTools` (Heusser et al., 2018b). Specifically, we used the `CountVectorizer` class to transform
646 the text from each window into a vector of word counts (using the union of all words across all
647 scenes as the “vocabulary,” excluding English stop words); this yielded a number-of-windows
648 by number-of-words *word count* matrix. We then used the `LatentDirichletAllocation` class
649 (`topics=100`, `method='batch'`) to fit a topic model (Blei et al., 2003) to the word count matrix,

yielding a number-of-windows (1047) by number-of-topics (100) *topic proportions* matrix. The topic proportions matrix describes the gradually evolving mix of topics (latent themes) present in each scene. Next, we transformed the topic proportions matrix to match the 1976 fMRI volume acquisition times. We assigned each topic vector to the timepoint (in seconds) midway between the beginning of the first scene and the end of the last scene in its corresponding sliding text window. By doing so, we warped the linear temporal distance between consecutive topic vectors to align with the inconsistent temporal distance between consecutive annotations (whose durations varied greatly). We then rescaled these timepoints to 1.5s TR units, and used linear interpolation to estimate a topic vector for each TR. This resulted in a number-of-TRs (1976) by number-of-topics (100) matrix.

We created similar topic proportions matrices using hand-annotated transcripts of each participant's verbal recall of the video (annotated by Chen et al., 2017). We tokenized the transcript into a list of sentences, and then re-organized the list into overlapping sliding windows spanning (up to) 10 sentences each, analogously to how we parsed the video annotations. In turn, we transformed each window's sentences into a word count vector (using the same vocabulary as for the video model), and then we used the topic model already trained on the video scenes to compute the most probable topic proportions for each sliding window. This yielded a number-of-windows (range: 83–312) by number-of-topics (100) topic proportions matrix for each participant. These reflected the dynamic content of each participant's recalls. Note: for details on how we selected the video and recall window lengths and number of topics, see *Supporting Information* and Figure S1.

Parsing topic trajectories into events using Hidden Markov Models

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

677 to implement this segmentation.

678 We used an optimization procedure to select the appropriate K for each topic proportions
679 matrix. Prior studies on narrative structure and processing have shown that we both perceive
680 and internally represent the world around us at multiple, hierarchical timescales (e.g., Hasson
681 et al., 2008; Lerner et al., 2011; Hasson et al., 2015; Chen et al., 2017; Baldassano et al., 2017, 2018).
682 However, for the purposes of our framework, we sought to identify the single timeseries of event-
683 representations that is emphasized *most heavily* in the temporal structure of the video and of each
684 participant's recall. We quantified this as the set of K states that maximized the similarity between
685 topic vectors for timepoints comprising each state, while minimizing the similarity between topic
686 vectors for timepoints across different states. Specifically, we computed (for each matrix)

$$\operatorname{argmax}_K [W_1(a, b)],$$

687 where a was the distribution of within-state topic vector correlations, and b was the distribution of
688 across-state topic vector correlations . We computed the first Wasserstein distance (W_1 ; also known
689 as *Earth mover's distance*; Dobrushin, 1970; Ramdas et al., 2017) between these distributions for a
690 large range of possible K -values (range [2, 50]), and selected the K that yielded the maximum value.
691 Figure 2B displays the event boundaries returned for the video, and Figure S4 displays the event
692 boundaries returned for each participant's recalls. See Figure S6 for the optimization functions
693 for the video and recalls. After obtaining these event boundaries, we created stable estimates of
694 the content represented in each event by averaging the topic vectors across timepoints between
695 each pair of event boundaries. This yielded a number-of-events by number-of-topics matrix for
696 the video and recalls from each participant.

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

698 In traditional list-learning experiments, participants view a list of items (e.g., words) and then recall
699 the items later. Our video-recall event matching approach affords us the ability to analyze memory
700 in a similar way. The video and recall events can be treated analogously to studied and recalled

701 “items” in a list-learning study. We can then extend classic analyses of memory performance and
702 dynamics (originally designed for list-learning experiments) to the more naturalistic video recall
703 task used in this study.

704 Perhaps the simplest and most widely used measure of memory performance is *accuracy*—i.e.,
705 the proportion of studied (experienced) items (in this case, video events) that the participant later
706 remembered. Chen et al. (2017) used this method to rate each participant’s memory quality by
707 computing the proportion of (50, manually identified) scenes mentioned in their recall. We found a
708 strong across-participants correlation between these independent ratings and the proportion of 30
709 HMM-identified video events matched to participants’ recalls (Pearson’s $r(15) = 0.71, p = 0.002$).
710 We further considered a number of more nuanced memory performance measures that are typically
711 associated with list-learning studies. We also provide a software package, Quail, for carrying out
712 these analyses (Heusser et al., 2017).

713 **Probability of first recall (PFR).** PFR curves (Welch and Burnett, 1924; Postman and Phillips,
714 1965; Atkinson and Shiffrin, 1968) reflect the probability that an item will be recalled first as a
715 function of its serial position during encoding. To carry out this analysis, we initialized a number-
716 of-participants (17) by number-of-video-events (30) matrix of zeros. Then for each participant, we
717 found the index of the video event that was recalled first (i.e., the video event whose topic vector
718 was most strongly correlated with that of the first recall event) and filled in that index in the matrix
719 with a 1. Finally, we averaged over the rows of the matrix, resulting in a 1 by 30 array representing
720 the proportion of participants that recalled an event first, as a function of the order of the event’s
721 appearance in the video (Fig. 3A).

722 **Lag conditional probability curve (lag-CRP).** The lag-CRP curve (Kahana, 1996) reflects the
723 probability of recalling a given item after the just-recalled item, as a function of their relative
724 encoding positions (or *lag*). In other words, a lag of 1 indicates that a recalled item was presented
725 immediately after the previously recalled item, and a lag of -3 indicates that a recalled item came 3
726 items before the previously recalled item. For each recall transition (following the first recall), we

727 computed the lag between the current recall event and the next recall event, normalizing by the
728 total number of possible transitions. This yielded a number-of-participants (17) by number-of-lags
729 (-29 to +29; 58 lags total excluding lags of 0) matrix. We averaged over the rows of this matrix to
730 obtain a group-averaged lag-CRP curve (Fig. 3B).

731 **Serial position curve (SPC).** SPCs (Murdock, 1962) reflect the proportion of participants that
732 remember each item as a function of the items' serial positions during encoding. We initialized
733 a number-of-participants (17) by number-of-video-events (30) matrix of zeros. Then, for each
734 recalled event, for each participant, we found the index of the video event that the recalled event
735 most closely matched (via the correlation between the events' topic vectors) and entered a 1 into
736 that position in the matrix. This resulted in a matrix whose entries indicated whether or not each
737 event was recalled by each participant (depending on whether the corresponding entires were
738 set to one or zero). Finally, we averaged over the rows of the matrix to yield a 1 by 30 array
739 representing the proportion of participants that recalled each event as a function of the events'
740 order appearance in the video (Fig. 3C).

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

750 **Semantic clustering scores.** Semantic clustering describes a participant's tendency to recall se-
751 mantically similar presented items together in their recall sequences. Here, we used the topic
752 vectors for each event as a proxy for its semantic content. Thus, the similarity between the seman-

tic content for two events can be computed by correlating their respective topic vectors. For each recall event transition, we sorted all not-yet-recalled events according to how correlated the topic vector of the closest-matching video event was to the topic vector of the closest-matching video event to the just-recalled event. We then computed the percentile rank of the observed next recall. We averaged these percentile ranks across all of the participant's recalls to obtain a single semantic clustering score for the participant.

759 Novel naturalistic memory metrics

760 **Precision.** We tested whether participants who recalled more events were also more *precise* in
761 their recollections. For each participant, we computed the average correlation between the topic
762 vectors for each recall event and those of its closest-matching video event. This gave a single value
763 per participant representing the average precision across all recalled events. We then correlated
764 these values with both hand-annotated and model-derived (i.e., the number of unique video events
765 matched by a participant's recall events) memory performance.

766 **Distinctiveness.** We also considered the *distinctiveness* of each recalled event. That is, how unique
767 a participant's description of a video event was, versus their descriptions of other video events.
768 We hypothesized that participants with high memory performance might describe each event in
769 a more distinctive way (relative to those with lower memory performance who might describe
770 events in a more general way). To test this hypothesis we define a distinctiveness score for each
771 recall event i as

$$d(i) = 1 - \frac{1}{N-1} \sum_{j=i} \text{corr}(\text{event}_i, \text{event}_j)$$

772 where the average is taken over the correlation between the recall event i 's topic vector and the
773 topic vectors from all other recall events from that participant. We averaged these distinctiveness
774 scores across all of the events recalled by the given participant to get the participant's distinctiveness
775 score. We correlated these distinctiveness scores with hand-annotated and model-derived memory

776 performance scores across-subjects, as above.

777 **Averaging correlations** In all instances where we performed statistical tests involving precision
778 or distinctiveness scores, we used the Fisher z -transformation (Fisher, 1925) to stabilize the variance
779 across the distribution of correlation values prior to performing the test. Similarly, when averaging
780 precision or distinctiveness scores, we z -transformed the scores prior to computing the mean, and
781 inverse z -transformed the result.

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

783 We used the UMAP algorithm (McInnes et al., 2018) to project the 100-dimensional topic space onto
784 a two-dimensional space for visualization (Figs. 7, 8). To ensure that all of the trajectories were
785 projected onto the *same* lower dimensional space, we computed the low-dimensional embedding
786 on a “stacked” matrix created by vertically concatenating the events-by-topics topic proportions
787 matrices for the video, across-participants average recall and all 17 individual participants’ recalls.
788 We then separated the rows of the result (a total-number-of-events by two matrix) back into
789 individual matrices for the video topic trajectory, across-participant average recall trajectory and the
790 trajectories for each individual participant’s recalls (Fig. 7). This general approach for discovering
791 a shared low-dimensional embedding for a collections of high-dimensional observations follows
792 Heusser et al. (2018b).

793 We optimized the manifold space for visualization based on two criteria: First, that the 2D
794 embedding of the video trajectory should reflect its original 100-dimensional structure as faithfully
795 as possible. Second, that the path traversed by the embedded video trajectory should intersect
796 itself a minimal number of times. The first criteria helps bolster the validity of visual intuitions
797 about relationships between sections of video content, based on their locations in the embedding
798 space. The second criteria was motivated by the observed low off-diagonal values in the video
799 trajectory’s temporal correlation matrix (suggesting that the same topic-space coordinates should
800 not be revisited; see Figure 2A in the main text). For further details on how we created this
801 low-dimensional embedding space, see *Supporting Information*.

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

803 In Figure 7B, we present an analysis aimed at characterizing locations in topic space that dif-
804 ferent participants move through in a consistent way (via their recall topic trajectories). The
805 two-dimensional topic space used in our visualizations (Fig. 7) comprised a 60 x 60 (arbitrary
806 units) square. We tiled this space with a 50 x 50 grid of evenly spaced vertices, and defined a
807 circular area centered on each vertex whose radius was two times the distance between adjacent
808 vertices (i.e., 2.4 units). For each vertex, we examined the set of line segments formed by connecting
809 each pair successively recalled events, across all participants, that passed through this circle. We
810 computed the distribution of angles formed by those segments and the x -axis, and used a Rayleigh
811 test to determine whether the distribution of angles was reliably “peaked” (i.e., consistent across
812 all transitions that passed through that local portion of topic space). To create Figure 7B we drew
813 an arrow originating from each grid vertex, pointing in the direction of the average angle formed
814 by the line segments that passed within 2.4 units. We set the arrow lengths to be inversely propor-
815 tional to the p -values of the Rayleigh tests at each vertex. Specifically, for each vertex we converted
816 all of the angles of segments that passed within 2.4 units to unit vectors, and we set the arrow
817 lengths at each vertex proportional to the length of the (circular) mean vector. We also indicated
818 any significant results ($p < 0.05$, corrected using the Benjamani-Hochberg procedure) by coloring
819 the arrows in blue (darker blue denotes a lower p -value, i.e., a longer mean vector); all tests with
820 $p \geq 0.05$ are displayed in gray and given a lower opacity value.

821 **Searchlight fMRI analyses**

822 In Figure 9, we present two analyses aimed at identifying brain regions whose responses (as par-
823 ticipants viewed the video) exhibited a particular temporal structure. We developed a searchlight
824 analysis wherein we constructed a 5 x 5 x 5 cube of voxels (following Chen et al., 2017) centered on
825 each voxel in the brain, and for each of these cubes, computed the temporal correlation matrix of
826 the voxel responses during video viewing. Specifically, for each of the 1976 volumes collected dur-
827 ing video viewing, we correlated the activity patterns in the given cube with the activity patterns

828 (in the same cube) collected during every other timepoint. This yielded a 1976 by 1976 correlation
829 matrix for each cube. Note: participant 5's scan ended 75s early, and in Chen et al., 2017's publicly
830 released dataset, their scan data was padded to match the length of the other participants'. For
831 our searchlight analyses, we removed this padded data (i.e., the last 50 TRs), resulting in a 1925 by
832 1925 correlation matrix for each cube in participant 5's brain.

833 Next, we constructed a series of "template" matrices. The first template reflected the timecourse
834 of the video's topic trajectory, and the others reflected the timecourse of each participant's recall
835 trajectory. To construct the video template, we computed the correlations between the topic
836 proportions estimated for every pair of TRs (prior to segmenting the trajectory into discrete events;
837 i.e., the correlation matrix shown in Figs. 2B and 9A). We constructed similar temporal correlation
838 matrices for each participant's recall topic trajectory (Figs. 2D, S4). However, to correct for length
839 differences and potential non-linear transformations between viewing time and recall time, we
840 first used dynamic time warping (Berndt and Clifford, 1994) to temporally align participants'
841 recall topic trajectories with the video topic trajectory. An example correlation matrix before and
842 after warping is shown in Fig. 9B. This yielded a 1976 by 1976 correlation matrix for the video
843 template and for each participant's recall template.

844 The temporal structure of the video's content (as described by our model) is captured in the
845 block-diagonal structure of the video's temporal correlation matrix (e.g., Figs. 2B, 9A), with time
846 periods of thematic stability represented as dark blocks of varying sizes. Inspecting the video
847 correlation matrix suggests that the video's semantic content is highly temporally specific (i.e., the
848 correlations between topic vectors from distant timepoints are almost all near zero). By contrast,
849 the activity patterns of individual (cubes of) voxels can encode relatively limited information
850 on their own, and their activity frequently contributes to multiple separate functions (Freedman
851 et al., 2001; Sigman and Dehaene, 2008; Charron and Koechlin, 2010; Rishel et al., 2013). By
852 nature, these two attributes give rise to similarities in activity across large timescales that may not
853 necessarily reflect a single task. To enable a more sensitive analysis of brain regions whose shifts
854 in activity patterns mirrored shifts in the semantic content of the video or recalls, we restricted
855 the temporal correlations we considered to the timescale of semantic information captured by our

model. Specifically, we isolated the upper triangle of the video correlation matrix and created a “proximal correlation mask” that included only diagonals from the upper triangle of the video correlation matrix up to the first diagonal that contained no positive correlations. Applying this mask to the full video correlation matrix was analogous to excluding diagonals beyond the corner of the largest diagonal block. In other words, the timescale of temporal correlations we considered corresponded to the longest period of thematic stability in the video, and by extension the longest expected period of thematic stability in participants’ recalls and the longest period of stability we might expect to see in voxel activity arising from processing or encoding video content. Figure 9 shows this proximal correlation mask applied to the temporal correlation matrices for the video, an example participant’s (warped) recall, and an example cube of voxels from our searchlight analyses.

To determine which (cubes of) voxel responses matched the video template, we correlated the proximal diagonals from the upper triangle of the voxel correlation matrix for each cube with the proximal diagonals from video template matrix (Kriegeskorte et al., 2008). This yielded, for each participant, a voxelwise map of correlation values. We then performed a one-sample *t*-test on the distribution of (Fisher z-transformed) correlations at each voxel, across participants. This resulted in a value for each voxel (cube), describing how reliably its timecourse followed that of the video.

We further sought to ensure that our analysis identified regions where the activations’ temporal structure specifically reflected that of the video, rather than regions whose activity was simply autocorrelated at a width similar to the video template’s diagonal. To achieve this, we used a phase shift-based permutation procedure, whereby we circularly shifted the video’s topic trajectory by a random number of timepoints, computed the resulting “null” video template, and re-ran the searchlight analysis, in full. (For each of the 100 permutations, the same random shift was used for all participants). We z-scored the observed (unshifted) result at each voxel against the distribution of permutation-derived “null” results, and estimated a *p*-value by computing the proportion of shifted results that yielded larger values. To create the map in Figure 9C, we thresholded out any voxels whose similarity to the unshifted video’s structure fell below the 95th percentile of the permutation-derived similarity results.

884 We used an analogous procedure to identify which voxels' responses reflected the recall tem-
885 plates. For each participant, we correlated the proximal diagonals from the upper triangle of the
886 correlation matrix for each cube of voxels with the proximal diagonals from the upper triangle
887 of their (time-warped) recall correlation matrix. As in the video template analysis, this yielded a
888 voxelwise map of correlation coefficients per participant. However, whereas the video analysis
889 compared every participant's responses to the same template, here the recall templates were unique
890 for each participant. As in the analysis described above, we *t*-scored the (Fisher *z*-transformed)
891 voxelwise correlations, and used the same permutation procedure we developed for the video
892 responses to ensure specificity to the recall timeseries and assign significance values. To create the
893 map in Figure 9D we again thresholded out any voxels whose scores were below the 95th percentile
894 of the permutation-derived null distribution.

895 **Neurosynth decoding analyses**

896 Neurosynth parses a massive online database of over 14,000 neuroimaging studies and constructs
897 meta-analysis images for over 13,000 psychology- and neuroscience-related terms, based on NIfTI
898 images accompanying studies where those terms appear at a high frequency. Given a novel image
899 (tagged with its value type; e.g., *t*-, *F*- or *p*-statistics), Neurosynth returns a list of terms whose
900 meta-analysis images are most similar. Our permutation procedure yielded, for each of the two
901 searchlight analyses, a voxelwise map of significance (*p*-statistic) values. These maps describe the
902 extent to which each voxel *specifically* reflected the temporal structure of the video or individuals'
903 recalls (i.e., for each voxel, the proportion of phase-shifted topic vector correlation matrices less
904 similar to the voxel activity correlation matrix than the unshifted video's correlation matrix). We
905 inputted the two statistical maps described above to Neurosynth to create a list of the 10 most
906 representative terms for each map.

907 **References**

- 908 Atkinson, R. C. and Shiffrin, R. M. (1968). Human memory: A proposed system and its control
909 processes. In Spence, K. W. and Spence, J. T., editors, *The psychology of learning and motivation*,
910 volume 2, pages 89–105. Academic Press, New York.
- 911 Baldassano, C., Chen, J., Zadbood, A., Pillow, J. W., Hasson, U., and Norman, K. A. (2017).
912 Discovering event structure in continuous narrative perception and memory. *Neuron*, 95(3):709–
913 721.
- 914 Baldassano, C., Hasson, U., and Norman, K. A. (2018). Representation of real-world event schemas
915 during narrative perception. *Journal of Neuroscience*, 38(45):9689–9699.
- 916 Berndt, D. J. and Clifford, J. (1994). Using dynamic time warping to find patterns in time series. In
917 *KDD workshop*, volume 10, pages 359–370.
- 918 Blei, D. M. and Lafferty, J. D. (2006). Dynamic topic models. In *Proceedings of the 23rd International
919 Conference on Machine Learning*, ICML '06, pages 113–120, New York, NY, US. ACM.
- 920 Blei, D. M., Ng, A. Y., and Jordan, M. I. (2003). Latent dirichlet allocation. *Journal of Machine
921 Learning Research*, 3:993 – 1022.
- 922 Brown, T. B., Mann, B., Ryder, N., Subbiah, M., Kaplan, J., Dhariwal, P., Neelakantan, A., Shyam,
923 P., Sastry, G., Askell, A., Agarwal, S., Herbert-Voss, A., Krueger, G., Henighan, T., Child, R.,
924 Ramesh, A., Ziegler, D. M., Wu, J., Winter, C., Hesse, C., Chen, M., Sigler, E., Litwin, M., Gray, S.,
925 Chess, B., Clark, J., Berner, C., McCandlish, S., Radford, A., Sutskever, I., and Amodei, D. (2020).
926 Language models are few-shot learners. *arXiv*, 2005.14165.
- 927 Brunec, I. K., Moscovitch, M. M., and Barene, M. D. (2018). Boundaries shape cognitive represen-
928 tations of spaces and events. *Trends in Cognitive Sciences*, 22(7):637–650.
- 929 Bylinskii, Z., Isola, P., Bainbridge, C., Torralba, A., and Oliva, A. (2015). Intrinsic and extrinsic
930 effects on image memorability. *Vision Research*, 116:165–178.

- 931 Capota, M., Turek, J., Chen, P.-H., Zhu, X., Manning, J. R., Sundaram, N., Keller, B., Wang, Y., and
932 Shin, Y. S. (2017). Brain imaging analysis kit.
- 933 Cer, D., Yang, Y., Kong, S. Y., Hua, N., Limtiaco, N., John, R. S., Constant, N., Guajardo-Cespedes,
934 M., Yuan, S., Tar, C., Sung, Y.-H., Strope, B., and Kurzweil, R. (2018). Universal sentence encoder.
935 *arXiv*, 1803.11175.
- 936 Charron, S. and Koechlin, E. (2010). Divided representations of current goals in the human frontal
937 lobes. *Science*, 328(5976):360–363.
- 938 Chen, J., Leong, Y. C., Honey, C. J., Yong, C. H., Norman, K. A., and Hasson, U. (2017). Shared
939 memories reveal shared structure in neural activity across individuals. *Nature Neuroscience*,
940 20(1):115.
- 941 Chun, M. and Turk-Browne, N. (2007). Interactions between attention and memory. *Current opinion
942 in neurobiology*, 17(2):177–184.
- 943 Clewett, D. and Davachi, L. (2017). The ebb and flow of experience determines the temporal
944 structure of memory. *Curr Opin Behav Sci*, 17:186–193.
- 945 Cohn-Sheely, B. I. and Ranganath, C. (2017). Time regained: how the human brain constructs
946 memory for time. *Current Opinion in Behavioral Sciences*, 17:169–177.
- 947 Dobrushin, R. L. (1970). Prescribing a system of random variables by conditional distributions.
948 *Theory of Probability & Its Applications*, 15(3):458–486.
- 949 DuBrow, S. and Davachi, L. (2013). The influence of contextual boundaries on memory for the
950 sequential order of events. *Journal of Experimental Psychology: General*, 142(4):1277–1286.
- 951 Ezzyat, Y. and Davachi, L. (2011). What constitutes an episode in episodic memory? *Psychological
952 Science*, 22(2):243–252.
- 953 Fisher, R. A. (1925). *Statistical Methods for Research Workers*. Oliver and Boyd.

- 954 Freedman, D., Riesenhuber, M., Poggio, T., and Miller, E. (2001). Categorical representation of
955 visual stimuli in the primate prefrontal cortex. *Science*, 291(5502):312–316.
- 956 Friendly, M., Franklin, P. E., Hoffman, D., and Rubin, D. C. (1982). The Toronto Word Pool:
957 Norms for imagery, concreteness, orthographic variables, and grammatical usage for 1,080
958 words. *Behavior Research Methods and Instrumentation*, 14:375–399.
- 959 Hasson, U., Chen, J., and Honey, C. J. (2015). Hierarchical process memory: memory as an integral
960 component of information processing. *Trends in Cognitive Science*, 19(6):304–315.
- 961 Hasson, U., Yang, E., Vallines, I., Heeger, D. J., and Rubin, N. (2008). A hierarchy of temporal
962 receptive windows in human cortex. *Journal of Neuroscience*, 28(10):2539–2550.
- 963 Heusser, A. C., Ezzyat, Y., Shiff, I., and Davachi, L. (2018a). Perceptual boundaries cause mnemonic
964 trade-offs between local boundary processing and across-trial associative binding. *Journal of*
965 *Experimental Psychology Learning, Memory, and Cognition*, 44(7):1075–1090.
- 966 Heusser, A. C., Fitzpatrick, P. C., Field, C. E., Ziman, K., and Manning, J. R. (2017). Quail: a
967 Python toolbox for analyzing and plotting free recall data. *The Journal of Open Source Software*,
968 10.21105/joss.00424.
- 969 Heusser, A. C., Ziman, K., Owen, L. L. W., and Manning, J. R. (2018b). HyperTools: a Python
970 toolbox for gaining geometric insights into high-dimensional data. *Journal of Machine Learning*
971 *Research*, 18(152):1–6.
- 972 Howard, M. W. and Kahana, M. J. (2002). A distributed representation of temporal context. *Journal*
973 *of Mathematical Psychology*, 46:269–299.
- 974 Howard, M. W., MacDonald, C. J., Tiganj, Z., Shankar, K. H., Du, Q., Hasselmo, M. E., and H., E.
975 (2014). A unified mathematical framework for coding time, space, and sequences in the medial
976 temporal lobe. *Journal of Neuroscience*, 34(13):4692–4707.
- 977 Howard, M. W., Viskontas, I. V., Shankar, K. H., and Fried, I. (2012). Ensembles of human MTL
978 neurons “jump back in time” in response to a repeated stimulus. *Hippocampus*, 22:1833–1847.

- 979 Huk, A., Bonnen, K., and He, B. J. (2018). Beyond trial-based paradigms: continuous behavior, on-
980 going neural activity, and naturalistic stimuli. *Journal of Neuroscience*, 10.1523/JNEUROSCI.1920-
981 17.2018.
- 982 Kahana, M. J. (1996). Associative retrieval processes in free recall. *Memory & Cognition*, 24:103–109.
- 983 Kelly, S., Lloyd, D., Nurmikko, T., and Roberts, N. (2007). Retrieving autobiographical memories
984 of painful events activates the anterior cingulate cortex and inferior frontal gyrus. *The Journal of
985 Pain*, 8(4):307–314.
- 986 Kriat, A. and Goldsmith, M. (1994). Memory in naturalistic and laboratory contexts: distin-
987 guishing accuracy-oriented and quantity-oriented approaches to memory assessment. *Journal of
988 Experimental Psychology: General*, 123(3):297–315.
- 989 Kriegeskorte, N., Mur, M., and Bandettini, P. (2008). Representational similarity analysis – con-
990 nnecting the branches of systems neuroscience. *Frontiers in Systems Neuroscience*, 2:1 – 28.
- 991 Landauer, T. K. and Dumais, S. T. (1997). A solution to Plato’s problem: the latent semantic
992 analysis theory of acquisition, induction, and representation of knowledge. *Psychological Review*,
993 104:211–240.
- 994 Lerner, Y., Honey, C. J., Silbert, L. J., and Hasson, U. (2011). Topographic mapping of a hierarchy
995 of temporal receptive windows using a narrated story. *Journal of Neuroscience*, 31(8):2906–2915.
- 996 Manning, J. R. (2019). Episodic memory: mental time travel or a quantum ‘memory wave’ function?
997 *PsyArXiv*, doi:10.31234/osf.io/6zjwb.
- 998 Manning, J. R. (2020). Context reinstatement. In Kahana, M. J. and Wagner, A. D., editors, *Handbook
999 of Human Memory*. Oxford University Press.
- 1000 Manning, J. R., Norman, K. A., and Kahana, M. J. (2015). The role of context in episodic memory.
1001 In Gazzaniga, M., editor, *The Cognitive Neurosciences, Fifth edition*, pages 557–566. MIT Press.

- 1002 Manning, J. R., Polyn, S. M., Baltuch, G., Litt, B., and Kahana, M. J. (2011). Oscillatory patterns
1003 in temporal lobe reveal context reinstatement during memory search. *Proceedings of the National*
1004 *Academy of Sciences, USA*, 108(31):12893–12897.
- 1005 McInnes, L., Healy, J., and Melville, J. (2018). UMAP: Uniform manifold approximation and
1006 projection for dimension reduction. *arXiv*, 1802(03426).
- 1007 Medford, N. and Critchley, H. D. (2010). Conjoint activity of anterior insular and anterior cingulate
1008 cortex: awareness and response. *Brain Structure and Function*, 214(5-6):535–549.
- 1009 Menon, V. and Uddin, L. Q. (2010). Saliency, switching, attention and control: a network model of
1010 insula function. *Brain Structure and Function*, 214(5-6):655–667.
- 1011 Mikolov, T., Chen, K., Corrado, G., and Dean, J. (2013). Efficient estimation of word representations
1012 in vector space. *arXiv*, 1301.3781.
- 1013 Mueller, A., Fillion-Robin, J.-C., Boidol, R., Tian, F., Nechifor, P., yoonsubKim, Peter, Rampin, R.,
1014 Corvellec, M., Medina, J., Dai, Y., Petrushev, B., Langner, K. M., Hong, Alessio, Ozsvald, I.,
1015 vkolmakov, Jones, T., Bailey, E., Rho, V., IgorAPM, Roy, D., May, C., foobuzz, Piyush, Seong,
1016 L. K., Goey, J. V., Smith, J. S., Gus, and Mai, F. (2018). WordCloud 1.5.0: a little word cloud
1017 generator in Python. *Zenodo*, <https://zenodo.org/record/1322068#.W4tPKZNKh24>.
- 1018 Murdock, B. B. (1962). The serial position effect of free recall. *Journal of Experimental Psychology*,
1019 64:482–488.
- 1020 Paller, K. A. and Wagner, A. D. (2002). Observing the transformation of experience into memory.
1021 *Trends in Cognitive Sciences*, 6(2):93–102.
- 1022 Pedregosa, F., Varoquaux, G., Gramfort, A., Michel, V., Thirion, B., Grisel, O., Blondel, M., Pretten-
1023 hofer, P., Weiss, R., Dubourg, V., Vanderplas, J., Passos, A., Cournapeau, D., Brucher, M., Perrot,
1024 M., and Duchesnay, E. (2011). Scikit-learn: Machine learning in Python. *Journal of Machine*
1025 *Learning Research*, 12:2825–2830.

- 1026 Polyn, S. M., Norman, K. A., and Kahana, M. J. (2009). A context maintenance and retrieval model
1027 of organizational processes in free recall. *Psychological Review*, 116(1):129–156.
- 1028 Postman, L. and Phillips, L. W. (1965). Short-term temporal changes in free recall. *Quarterly Journal*
1029 of *Experimental Psychology*, 17:132–138.
- 1030 Rabiner, L. (1989). A tutorial on Hidden Markov Models and selected applications in speech
1031 recognition. *Proceedings of the IEEE*, 77(2):257–286.
- 1032 Radford, A., Wu, J., Child, R., Luan, D., Amodei, D., and Sutskever, I. (2019). Language models are
1033 unsupervised multitask learners. *OpenAI Blog*, 1(8).
- 1034 Radvansky, G. A. and Zacks, J. M. (2017). Event boundaries in memory and cognition. *Curr Opin*
1035 *Behav Sci*, 17:133–140.
- 1036 Ramdas, A., Trillos, N., and Cuturi, M. (2017). On wasserstein two-sample testing and related
1037 families of nonparametric tests. *Entropy*, 19(2):47.
- 1038 Ranganath, C. and Ritchey, M. (2012). Two cortical systems for memory-guided behavior. *Nature*
1039 *Reviews Neuroscience*, 13:713 – 726.
- 1040 Rishel, C. A., Huang, G., and Freedman, D. J. (2013). Independent category and spatial encoding
1041 in parietal cortex. *Neuron*, 77(5):969–979.
- 1042 Sigman, M. and Dehaene, S. (2008). Brain mechanisms of serial and parallel processing during
1043 dual-task performance. *Journal of Neuroscience*, 28(30):7585–7589.
- 1044 Simony, E., Honey, C. J., Chen, J., and Hasson, U. (2016). Uncovering stimulus-locked network
1045 dynamics during narrative comprehension. *Nature Communications*, 7(12141):1–13.
- 1046 Stark, S. M., Yassa, M. A., and Stark, C. E. L. (2010). Individual differences in spatial pattern
1047 separation performance associated with healthy aging in humans. *Learning & Memory*, 17(6):284–
1048 288.

- 1049 Tomrary, A. and Davachi, L. (2017). Consolidation promotes the emergence of representational
1050 overlap in the hippocampus and medial prefrontal cortex. *Neuron*, 96(1):228–241.
- 1051 Tyng, C. M., Amin, H. U., Saad, M. N. M., and S, M. A. (2017). The influences of emotion on
1052 learning and memory. *Frontiers in psychology*, 8:1454.
- 1053 Welch, G. B. and Burnett, C. T. (1924). Is primacy a factor in association-formation. *American Journal
1054 of Psychology*, 35:396–401.
- 1055 Yarkoni, T., Poldrack, R. A., Nichols, T. E., Van Essen, D. C., and Wager, T. D. (2011). Large-scale
1056 automated synthesis of human functional neuroimaging data. *Nature Methods*, 8(8):665.
- 1057 Yassa, M. A., Lacy, J. W., Stark, S. M., Albert, M. S., Gallagher, M., and Stark, C. E. L. (2011). Pattern
1058 separation deficits associated with increased hippocampal ca3 and dentate gyrus activity in
1059 nondemented older adults. *Hippocampus*, 21(9):968–979.
- 1060 Yassa, M. A. and Stark, C. E. L. (2011). Pattern separation in the hippocampus. *Trends In Neuro-
1061 sciences*, 34(10):515–525.
- 1062 Yonelinas, A. P. (2002). The nature of recollection and familiarity: A review of 30 years of research.
1063 *Journal of Memory and Language*, 46:441–517.
- 1064 Yonelinas, A. P., Kroll, N. E., Quamme, J. R., Lazzara, M. M., Sauvé, M. J., Widaman, K. F., and
1065 Knight, R. T. (2002). Effects of extensive temporal lobe damage or mild hypoxia on recollection
1066 and familiarity. *Nature Neuroscience*, 5(11):1236–41.
- 1067 Zacks, J. M., Speer, N. K., Swallow, K. M., Braver, T. S., and Reynolds, J. R. (2007). Event perception:
1068 a mind-brain perspective. *Psychological Bulletin*, 133:273–293.
- 1069 Zadbood, A., Chen, J., Leong, Y. C., Norman, K. A., and Hasson, U. (2017). How we transmit
1070 memories to other brains: Constructing shared neural representations via communication. *Cereb
1071 Cortex*, 27(10):4988–5000.
- 1072 Zwaan, R. A. and Radvansky, G. A. (1998). Situation models in language comprehension and
1073 memory. *Psychological Bulletin*, 123(2):162 – 185.

1074 **Supporting information**

1075 Supporting information is available in the online version of the paper.

1076 **Acknowledgements**

1077 We thank Luke Chang, Janice Chen, Chris Honey, Lucy Owen, Emily Whitaker, and Kirsten Ziman
1078 for feedback and scientific discussions. We also thank Janice Chen, Yuan Chang Leong, Kenneth
1079 Norman, and Uri Hasson for sharing the data used in our study. Our work was supported in part
1080 by NSF EPSCoR Award Number 1632738. The content is solely the responsibility of the authors
1081 and does not necessarily represent the official views of our supporting organizations.

1082 **Author contributions**

1083 Conceptualization: A.C.H. and J.R.M.; Methodology: A.C.H., P.C.F. and J.R.M.; Software: A.C.H.,
1084 P.C.F. and J.R.M.; Analysis: A.C.H., P.C.F. and J.R.M.; Writing, Reviewing, and Editing: A.C.H.,
1085 P.C.F. and J.R.M.; Supervision: J.R.M.

1086 **Author information**

1087 The authors declare no competing financial interests. Correspondence and requests for materials
1088 should be addressed to J.R.M. (jeremy.r.manning@dartmouth.edu).