

1 How is experience transformed into memory? Memory
2 for television episodes preserves event content while
3 introducing new across-event similarities

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 March 13, 2020

6 **Abstract**

7 The ways our Our experiences unfold over time define defining unique *trajectories* through
8 the relevant representational spaces. Within this geometric framework By casting our life events
9 as temporally evolving trajectories, one can compare the shape of the trajectory formed by an
10 experience to that defined by our later remembering of that experience. We propose a framework
11 for mapping naturalistic experiences onto geometric spaces that characterize how they unfold
12 experiences are segmented into discrete events, and how the contents of event sequences evolve
13 over time. We apply this approach to a naturalistic memory experiment which had participants
14 view and recount a video. We found that the shapes of the trajectories formed by television
15 episode. The content of participants' recounts were all highly similar to that of the original

16 video, but participants differed in the level of detail they remembered. We also of events from
17 the original episode closely matched the original episode's content. Further, we introduce two
18 novel metrics for assessing memory quality (precision and distinctiveness), both of which relate
19 to participants' ability to recapitulate the experience. Lastly, we identified a network of brain
20 structures that are sensitive to the "shapes" of our ongoing experiences, and an overlapping
21 network that is sensitive to how we will later remember those experiences. (at the time of
22 encoding) to how people later remembered those experiences in relation to other experiences. In
23 this way, modeling the content of richly structured experiences can reveal how (geometrically
24 and conceptually) those experiences are segmented into events and integrated into our memories
25 of other experiences.

26 Introduction

27 What does it mean to *remember* something? In traditional episodic memory experiments (e.g.,
28 list-learning or trial-based experiments; Murdock, 1962; Kahana, 1996), remembering is often cast
29 as a discrete and binary operation: each studied item may be separated from the rest of one's
30 experiences, and that item may be all others, and labeled as having been recalled versus or forgotten.
31 More nuanced studies might incorporate self-reported confidence measures as a proxy for memory
32 strength, or ask participants to discriminate between "recollecting" the (contextual) details of an
33 experience or having a general feeling of "familiarity" (Yonelinas, 2002). Using well-controlled,
34 trial-based experimental designs, the field has amassed a wealth of valuable information regarding
35 human episodic memory. However, characterizing and evaluating memory in more realistic
36 contexts (e.g., recounting a recent experience to a friend) is fundamentally different in at least
37 three ways there are fundamental properties of the external world and our memories that trial-based
38 experiments are not well suited to capture (for review also see Koriat and Goldsmith, 1994; Huk
39 et al., 2018). First, real-world recall is our experiences and memories are continuous, rather than
40 discrete. Unlike in trial-based experiments, removing discrete—removing a (naturalistic) event
41 from the context in which it occurs can substantially change its meaning. Second, the specific
42 words language used to describe an experience have has little bearing on whether the experi-

ence should be considered to have been “remembered.” Asking whether the rememberer has precisely reproduced a specific set of words to describe a given experience is nearly orthogonal to whether they were actually able to remember it. In classic (e.g., list-learning) memory studies, by contrast, ~~counting~~ the number or proportion of precise recalls is often a primary metric ~~of for~~ assessing the quality of participants’ memories. Third, one might remember the ~~gist or essence~~ ~~or a general summary~~ of an experience but forget (or neglect to recount) particular details. Capturing the ~~gist essence~~ of what happened is typically the main “point” of recounting a memory to a listener~~whereas, depending on the circumstances, accurate recall of any specific detail may be irrelevant. There is no analog of the gist of an experience in most traditional memory experiments; rather we tend to assess participants’ abilities to recover specific details (e.g., computing the proportion of specific stimuli they remember, which presentation positions the remembered stimuli came from, etc.).~~, while the addition of highly specific details may add comparatively little to ~~successful conveyance of an experience.~~

How might one go about formally characterizing the ~~gist~~ “essence” of an experience, or whether ~~that gist it~~ has been recovered by the rememberer? Any given moment of an experience derives meaning from surrounding moments, as well as from longer-range temporal associations (e.g., Lerner et al., 2011). Therefore (Lerner et al., 2011; Manning, 2019). Therefore, the time-course describing how an event unfolds is fundamental to its overall meaning. Further, this hierarchy formed by our subjective experiences at different timescales defines a *context* for each new moment (e.g., ?Howard et al., 2014)(e.g., Howard and Kahana, 2002; Howard et al., 2014), and plays an important role in how we interpret that moment and remember it later (for review see Manning et al., 2015). Our memory systems can ~~then~~ leverage these associations to form predictions that help guide our behaviors (Ranganath and Ritchey, 2012). For example, as we navigate the world, the features of our subjective experiences tend to change gradually (e.g., the room or situation we are in ~~at any given moment~~ is strongly temporally autocorrelated), allowing us to form stable estimates of our current situation and behave accordingly (Zacks et al., 2007; Zwaan and Radvansky, 1998). ~~Although our experiences most often change gradually, they also occasionally change suddenly~~

71 Occasionally, this gradual “drift” of our ongoing experience is punctuated by sudden changes,
72 or “shifts” (e.g., when we walk through a doorway; Radvansky and Zacks, 2017). Prior research
73 suggests that these sharp transitions (termed *event boundaries*) ~~during an experience~~ help to dis-
74cretize our experiences ~~(and their mental representations)~~ into *events* (Radvansky and Zacks, 2017;
75 Brunec et al., 2018; Heusser et al., 2018a; Clewett and Davachi, 2017; Ezzyat and Davachi, 2011;
76 DuBrow and Davachi, 2013). The interplay between the stable (~~within-event~~within-event) and
77 transient (~~across-event~~across-event) temporal dynamics of an experience also provides a potential
78 framework for transforming experiences into memories that distill those experiences down to their
79 ~~essence—i.e., their gists.~~ essence. For example, prior work has shown that event boundaries can
80 influence how we learn sequences of items (Heusser et al., 2018a; DuBrow and Davachi, 2013), nav-
81 igitate (Brunec et al., 2018), and remember and understand narratives (Zwaan and Radvansky, 1998;
82 Ezzyat and Davachi, 2011). Prior research has implicated a network of brain regions (including the
83 hippocampus and the medial prefrontal cortex) as playing a critical role in transforming experiences
84 into structured and consolidated memories (Tompry and Davachi, 2017).

85 Here we sought to examine how the temporal dynamics of a “naturalistic” experience were later
86 reflected in participants’ ~~later memories of that experience~~memories. We analyzed an open dataset
87 that comprised behavioral and functional Magnetic Resonance Imaging (fMRI) data collected as
88 participants viewed and then verbally ~~recalled~~recounted an episode of the BBC television series
89 *Sherlock* (Chen et al., 2017). We developed a computational framework for characterizing the
90 temporal dynamics of the moment-by-moment content of the episode, and of participants’ verbal
91 recalls. Specifically, we use topic modeling (Blei et al., 2003) to characterize the thematic conceptual
92 (semantic) content present in each moment of the episode and recalls, and ~~we use~~ Hidden Markov
93 Models (Rabiner, 1989; Baldassano et al., 2017) to discretize ~~the this~~ evolving semantic content
94 into events. In this way, we cast naturalistic experiences (and recalls of those experiences) as
95 geometric topic trajectories that describe how the experiences evolve over time. ~~In other words,~~
96 ~~the episode’s topic trajectory is a formalization of its gist.~~ Under this framework, successful
97 remembering entails verbally “traversing” the ~~topic content~~ trajectory of the ~~original~~ episode,
98 thereby reproducing the ~~original episode’s gist.~~ In addition, comparing shape (or essence) of the

99 original experience. Comparing the shapes of the topic trajectories of the original episode and of
100 participants' retellings of the episode then reveals which aspects of the episode were preserved (or
101 lost) in the translation into memory. We also identified a network of further introduce two novel
102 metrics for assessing memory quality: the precision with which a participant recounts each event
103 and 2) the distinctiveness of each recall event (relative to other recalled events). We examine how
104 these metrics relate to participants' overall memory performance, and discuss the ways in which
105 they improve upon classic "proportion-recalled" measures for analyzing naturalistic memory. Last,
106 we utilize our framework to identify networks of brain structures whose responses (as participants
107 watched the episode) reflected the gist temporal dynamics of the episode, and a second network
108 whose responses reflected how participants would later recount the episode it.

109 Results

110 To characterize the gists temporally dynamic contents of the *Sherlock* episode participants watched
111 and their subsequent recounts of the episode and participants' subsequent recounts, we used
112 a topic model (Blei et al., 2003) to discover the latent thematic content in the video themes. Topic
113 models take as inputs a vocabulary of words to consider and a collection of text documents; they
114 return as output two, and return two output matrices. The first output of these is a topics matrix
115 whose rows are topics (latent themes) and whose columns correspond to words in the vocabulary.
116 The entries of the topics matrix define how each word in the vocabulary is weighted by each
117 discovered topic. For example, a detective-themed topic might weight heavily on words like
118 "crime," and "search." The second output is a topic proportions matrix, with one row per document
119 and one column per topic. The topic proportions matrix describes which mix what mixture of
120 discovered topics is reflected in each document.

121 Chen et al. (2017) collected hand-annotated information about each of 1000 (manually identified)
122 scenes time segments spanning the roughly 45–50 minute video used in their experiment. This
123 information included: a brief narrative description of what was happening; whether the the
124 location where the scene took place indoors vs. outdoors; the names of any characters on the

125 screen; names of any characters who were in focus in the camera shot; names of characters who
126 were speaking; the location where the scene took place; the camera angle (close up, medium, long,
127 etc.); whether or not background music was present; and other similar details (for a full list of
128 annotated features, see *Methods*). We took from these annotations the union of all unique words
129 (excluding stop words, such as “and,” “or,” “but,” etc.) across all features and scenes as the
130 “vocabulary” for the topic model. We then concatenated the sets of words across all features
131 contained in overlapping 50-scene sliding windows, sliding windows of (up to) 50 scenes, and
132 treated each 50-scene sequence window as a single “document” for the purposes for the purpose
133 of fitting the topic model. Next, we fit a topic model with (up to) $K = 100$ topics to this collection of
134 documents. We found that 27–32 unique topics (with non-zero weights) were sufficient to describe
135 the time-varying content of the video (see *Methods*; Figs. 1, S2). Note that our approach is similar in
136 some respects to Dynamic Topic Models (Blei and Lafferty, 2006), in that we sought to characterize
137 how the thematic content of the episode evolved over time. However, whereas Dynamic Topic
138 Models are designed to characterize how the properties of *collections* of documents change over
139 time, our sliding window approach allows us to examine the topic dynamics within a single
140 document (or video). Specifically, our approach yielded (via the topic proportions matrix) a single
141 *topic vector* for each timepoint of the episode (we set timepoints sliding window of annotations
142 transformed by the topic model. We then stretched (interpolated) the resulting windows-by-topics
143 matrix to match the acquisition times time series of the 1976 fMRI volumes collected as participants
144 viewed the episode).

145 The 32 topics we found were heavily character-focused (e.g., the top-weighted word in
146 each topic was nearly always a character) and could be roughly divided into themes that were
147 primarily Sherlock Holmes focused (Sherlock is centered around Sherlock Holmes (the titular
148 character); primarily John Watson focused (John is, John Watson (Sherlock’s close confidant and
149 assistant); or that involved Sherlock and John interacting (supporting characters (e.g., Inspector
150 Lestrade, Sergeant Donovan, or Sherlock’s brother Mycroft), or the interactions between various
151 pairs of these characters (see Fig. S2). Several of the identified topics were highly similar, which we
152 hypothesized might allow us to distinguish between subtle narrative differences (if the distinctions

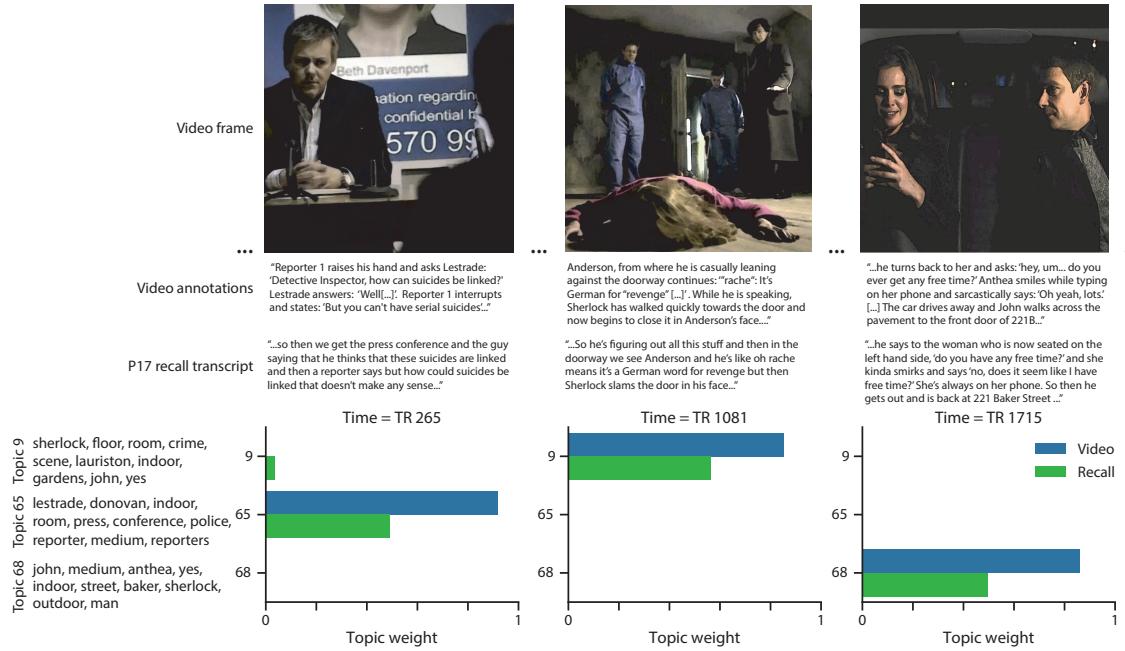


Figure 1: Methods overview. We used hand-annotated descriptions of each moment of video to fit a topic model. Three example video frames and their associated descriptions are displayed (top two rows). Participants later recalled the video (in the third row, we show example recalls of the same three scenes from participant [1317](#)). We used the topic model (fit to the annotations) to estimate topic vectors for each moment of video and each sentence the participants recalled. Example topic vectors are displayed in the bottom row (blue: video annotations; green: example participant’s recalls). Three topic dimensions are shown (the highest-weighted topics for each of the three example scenes, respectively). We also show the [ten](#) [10](#) highest-weighted words for each topic. Figure S2 provides a full list of the top 10 words from each of the discovered topics.

153 between those overlapping topics were meaningful; ~~also see Fig. S3~~. The topic vectors for each
154 timepoint were *sparse*, in that only a small number (usually one or two) of topics tended to be
155 “active” in any given timepoint (Fig. 2A). Further, the dynamics of the topic activations appeared
156 to exhibit ~~persistence~~persistence (i.e., given that a topic was active in one timepoint, it was likely to be
157 active in the following timepoint) along with *occasional rapid changes* (i.e., occasionally topics would
158 appear to spring into or out of existence). These two properties of the topic dynamics may be seen in
159 the block diagonal structure of the timepoint-by-timepoint correlation matrix (Fig. 2B). ~~Following~~
160 ~~Baldassano et al. (2017), we and reflect the gradual drift and sudden shifts fundamental to the~~
161 ~~temporal dynamics of real-world experiences. Given this observation, we adapted an approach~~
162 ~~devised by Baldassano et al. (2017), and~~ used a Hidden Markov Model (HMM) to identify the *event*
163 *boundaries* where the topic activations changed rapidly (i.e., at the boundaries of the blocks in the
164 correlation matrix; event boundaries identified by the HMM are outlined in yellow in Fig. 2B). Part
165 of our model fitting procedure required selecting an appropriate number of “events” ~~to segment~~
166 ~~the timeseries into. We~~ into which the topic trajectory should be segmented. To accomplish this,
167 ~~we~~ used an optimization procedure ~~to identify the number of events that maximized within-event~~
168 ~~stability while also minimizing across-event correlations that maximized the difference between the~~
169 ~~topic weights for timepoints within an event and across multiple events~~ (see *Methods* for additional
170 details). ~~To create~~ We then created a stable “summary” of the ~~video, we computed the average~~
171 ~~topic vector within each event~~ content within each video event by averaging the topic vectors
172 ~~across timepoints each event spanned~~ (Fig. 2C).

173 Given that the time-varying content of the video could be segmented cleanly into discrete
174 events, we wondered whether participants’ recalls of the video also displayed a similar structure.
175 We applied the same topic model (already trained on the video annotations) to each participant’s
176 recalls. Analogous to how we ~~analyzed~~parsed the time-varying content of the video, to obtain
177 similar estimates for ~~participants’ recalls~~each participant’s recall, we treated each ~~(overlapping)~~
178 ~~overlapping “window” of (up to 10 sentence “window” of) sentences from~~ their transcript as a
179 “document” ~~and then~~ and computed the most probable mix of topics reflected in each time-
180 point’s sentences. This yielded, for each participant, a ~~number of sentences~~number of windows

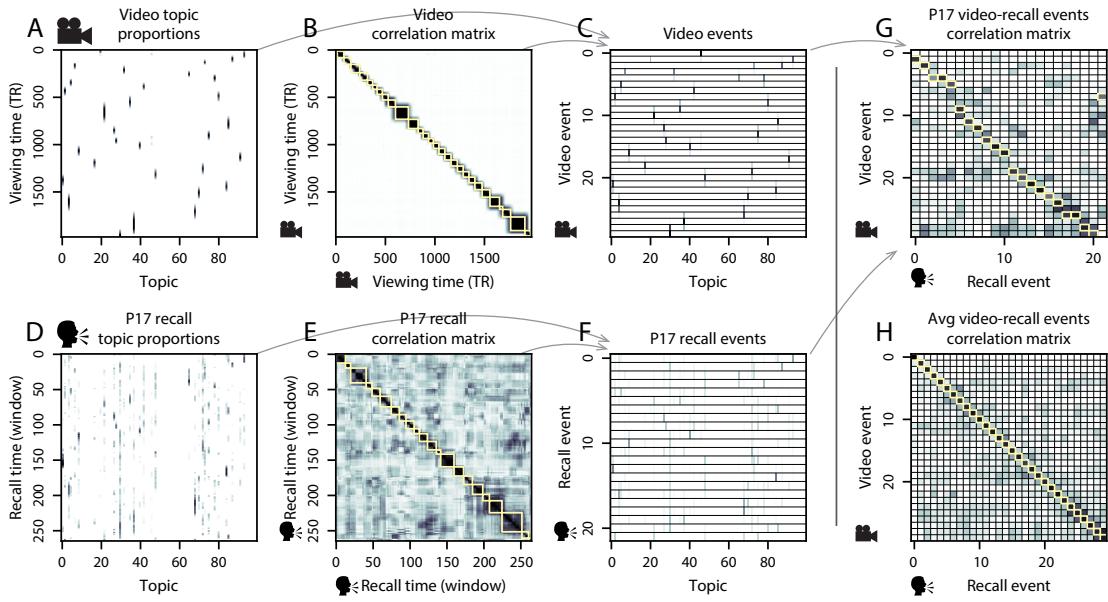


Figure 2: Modelling naturalistic stimuli and recalls. All panels: darker colors indicate greater values; range: [0, 1]. **A.** Topic vectors ($K = 100$) for each of the 1976 video timepoints. **B.** Timepoint-by-timepoint correlation matrix of the topic vectors displayed in Panel A. Event boundaries ~~detected~~-discovered by the HMM are denoted in yellow (34–30 events detected). **C.** Average topic vectors for each of the 34–30 video events. **D.** Topic vectors for each of 294–265 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 (27–22 events detected). For similar plots for all participants see Figure S4. **F.** Average topic vectors for each of the 27–22 recalled events from the example participant. **G.** Correlations between the topic vectors for every pair of video events (Panel C) and recalled events (from the example participant; Panel F). For similar plots for all participants, see Figure S5. **H.** Average correlations between each pair of video events and recalled events (across all 17 participants). To create the figure, each recalled event was assigned to the video event with the most correlated topic vector (yellow boxes in panels G and H). The heat maps in each panel were created using Seaborn (Waskom et al., 2016).

181 by number-of-topics topic proportions matrix that characterized how the topics identified in the
182 original video were reflected in the participant's recalls. Note that an important feature of our
183 approach is that it allows us to compare ~~participant's participants'~~ recalls to events from the orig-
184 inal video, despite ~~that different participants may have used different~~ different participants using
185 widely varying language to describe the same event, and that those descriptions may not match
186 the original annotations. This is a ~~huge~~ substantial benefit of projecting the video and recalls into
187 a shared "topic" space. An example topic proportions matrix from one participant's recalls is
188 shown in Figure 2D.

189 Although the example participant's recall topic proportions matrix has some visual similarity to
190 the video topic proportions matrix, the time-varying topic proportions for the example participant's
191 recalls are not as sparse as those for the video (e.g., compare Figs. 2A and D). Similarly, although
192 there do appear to be periods of stability in the recall topic dynamics (e.g.i.e., most topics are active
193 or inactive over contiguous blocks of time), the individual topics' overall timecourses are not as
194 cleanly delineated as the video topicsare. To examine these patterns in detail, we computed the
195 timepoint-by-timepoint correlation matrix for the example participant's recall topic ~~proportions~~
196 trajectory (Fig. 2E). As in the video correlation matrix (Fig. 2B), the example participant's recall
197 correlation matrix has a strong block diagonal structure, indicating that their recalls are discretized
198 into separated events. As for the video correlation matrix, we ~~can use an HMM, along with~~
199 ~~the aforementioned number-of-events~~ leveraged an HMM-based optimization procedure (also
200 see *Methods*) to determine how many events are reflected in the participant's recalls and where
201 specifically the event boundaries fall (outlined in yellow). We carried out a similar analysis on all
202 17 participants' recall topic proportions matrices (Fig. S4).

203 Two clear patterns emerged from this set of analyses. First, although every individual partic-
204 ipant's recalls could be segmented into discrete events (i.e., every individual participant's recall
205 correlation matrix exhibited clear block diagonal structure; Fig. S4), each participant appeared to
206 have a unique *recall resolution*, reflected in the sizes of those blocks. ~~For example~~ While, some par-
207 ticipants' recall topic proportions segmented into just a few events (e.g., Participants P1, P4, and
208 P15), ~~while others' recalls~~ P5, and P7, others' segmented into many shorter duration events (e.g.,

209 Participants P12, P13, and P17). This suggests that different participants may be recalling the video
210 with different levels of detail—e.g., some might touch on just the major plot points, whereas others
211 might attempt to recall every minor scene or action. The second clear pattern present in every in-
212 dividual participant’s recall correlation matrix is that, unlike in the video correlation matrix, there
213 are substantial off-diagonal correlations in participant’s recalls. Whereas each event in the original
214 video (~~was was~~) largely) separable from the others (Fig. 2B), in transforming those separable events
215 into memory, participants appear to be integrating ~~across different~~ across multiple events, blend-
216 ing elements of previously recalled and not-yet-recalled events content into each newly recalled
217 event (Figs. 2D, S4; also see Manning et al., 2011; Howard et al., 2012) (Figs. 2E, S4; also see Manning et al., 2011; Howar-
218 .

219 The above results indicate that both the structure of the original video and participants’ recalls
220 of the video exhibit event boundaries that can be identified automatically by characterizing the dy-
221 namic content using a shared topic model and segmenting the content into events using via HMMs.
222 Next, we asked whether some correspondence might be made between the specific content of the
223 events the participants experienced in the video, and the events they later recalled. One approach
224 to linking the experienced (video) and recalled events is to label each recalled event as matching the
225 video event with the most similar (i.e., most highly correlated) topic vector (Figs. 2G, S5). This yields
226 a sequence of “presented” events from the original video, and a sequence of (potentially differently
227 ordered) “recalled” events for each participant. Analogous to classic list-
228 learning studies, we can then examine participants’ recall sequences by asking which events they
229 tended to recall first (probability of first recall; Fig. A; Welch and Burnett, 1924; Postman and Phillips, 1965; Atkinson and
230 (probability of first recall; Fig. 3A; Atkinson and Shiffrin, 1968; Postman and Phillips, 1965; Welch and Burnett, 1924)
231 ; how participants most often transition between recalls of the events as a function of the temporal
232 distance between them (lag-conditional response probability; Fig. B; Kahana, 1996) (lag-conditional response probability
233 ; and which events they were likely to remember overall (serial position recall analyses; Fig. C; Murdoeck, 1962)
234 . In (serial position recall analyses; Fig. 3C; Murdock, 1962). Interestingly, for two of these analyses
235 (probability of first recall and lag-conditional response probability curves) we observe patterns
236 comparable to classic effects from the list-learning studies, this set of three analyses may be used to

237 gain a nearly complete view into the sequences of recalls participants made (e.g., Kahana, 2012).
238 Extending these analyses to apply to naturalistic stimuli and recall (Heusser et al., 2017) highlights
239 that, in naturalistic recall, these analyses provide a wholly incomplete picture: they leave out any
240 attempt to quantify participants' abilities to capture the *content* of what occurred in the video—their
241 only experimental instruction! literature: namely, a higher probability of initiating recall with the
242 first event in the sequence (Fig. 3A) and a higher probability of transitioning to neighboring events
243 with an asymmetric forward bias (Fig. 3B). In contrast, we do not observe a pattern comparable to
244 the serial position effect (Fig. 3C), but rather we see higher memory for specific events distributed
245 somewhat evenly throughout the video.

246 **The dynamic** We can also apply two list-learning-native analyses that describe how participants
247 group items in their recall sequences: temporal clustering and semantic clustering (Polyn et al., 2009, see *Methods* for deta
248 . Temporal clustering refers to the extent to which participants group their recall responses
249 according to encoding position. Overall, we found that sequentially viewed video events were
250 clustered heavily in participants' recall event sequences (mean clustering score: 0.767, SEM: 0.029),
251 and that participants with higher temporal clustering scores tended to perform better according
252 to both Chen et al. (2017)'s hand-annotated memory scores (Pearson's $r(15) = 0.62$, $p = 0.008$) and
253 our model's estimate (Pearson's $r(15) = 0.54$, $p = 0.024$). Semantic clustering measures the extent
254 to which participants cluster their recall responses according to semantic similarity. We found that
255 participants tended to recall semantically similar video events together (mean clustering score:
256 0.787, SEM: 0.018), and that semantic clustering score was also related to both hand-annotated
257 (Pearson's $r(15) = 0.65$, $p = 0.004$) and model-derived (Pearson's $r(15) = 0.63$, $p = 0.007$) memory
258 performance.

259 Statistical models of memory studies often treat recall success as binary (i.e., an item either was
260 or was not recalled), or occasionally categorical (e.g., to distinguish familiarity from recollection; Yonelinas et al., 2002)
261 . Such approaches are tenable in classical list-learning or recognition memory paradigms, as the
262 presented stimuli tend to be very simple (e.g., a sequence of individual words or items). However,
263 the feature-rich content of a naturalistic experiences may later be described with many, highly
264 variable levels of success. Our framework produces a content-based model of individual stimulus

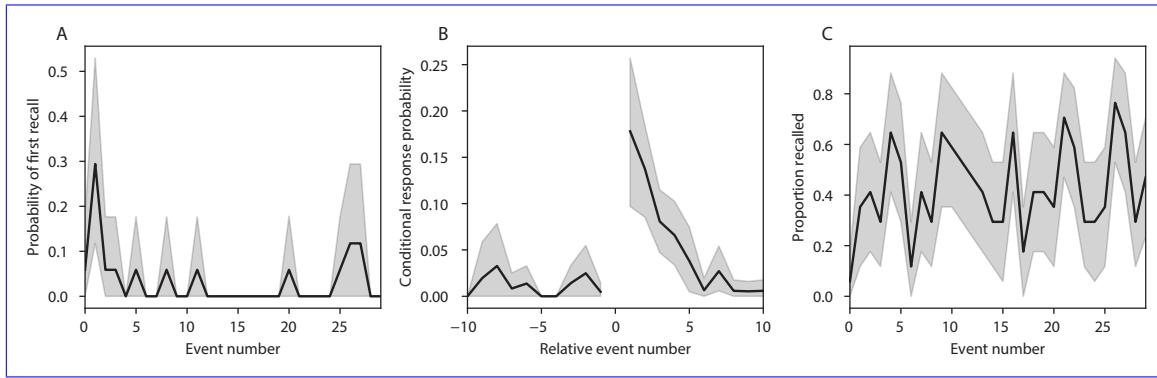


Figure 3: Naturalistic extensions of classic list-learning memory analyses. **A.** The probability of first recall as a function of the serial position of the event in the video. **B.** The probability of recalling each event, conditioned on having most recently recalled the event *lag* events away in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

and recall events by projecting the dynamic content of the video and participants' recalls into a shared topic space. This allows for direct, quantitative comparison between all stimulus and recall events, as well as between the recall events themselves. Leveraging these content-based models of the stimulus/recall events, we developed two novel, *continuous* metrics for analyzing naturalistic memory: *precision* and *distinctiveness*. We define precision as the “completeness” of recall, or how fully the presented content was recapitulated in memory. Under our framework, we quantify this for a given recall event as the correlation between the topic proportions of the recall event and the maximally correlated video event (Fig. 4). A second novel metric we introduce here is *distinctiveness*, which we define as the “specificity” of recall, or how unique the description of a given section of content was, compared to descriptions for other sections of content. We quantify this for each recall event as 1 minus the average correlation between the given recall event and all other recall events not matched to the same video event. In addition to individual events, one may also use these metrics to describe each participant's overall performance (i.e., by averaging across a participant's event-wise precision or distinctiveness scores). Participants whose recall events are more veridical descriptions of what happened in the video event will presumably have higher precision scores. We find that, across participants, a higher precision score is correlated to both hand-annotated memory performance (Pearson's $r(15) = 0.60, p = 0.010$)

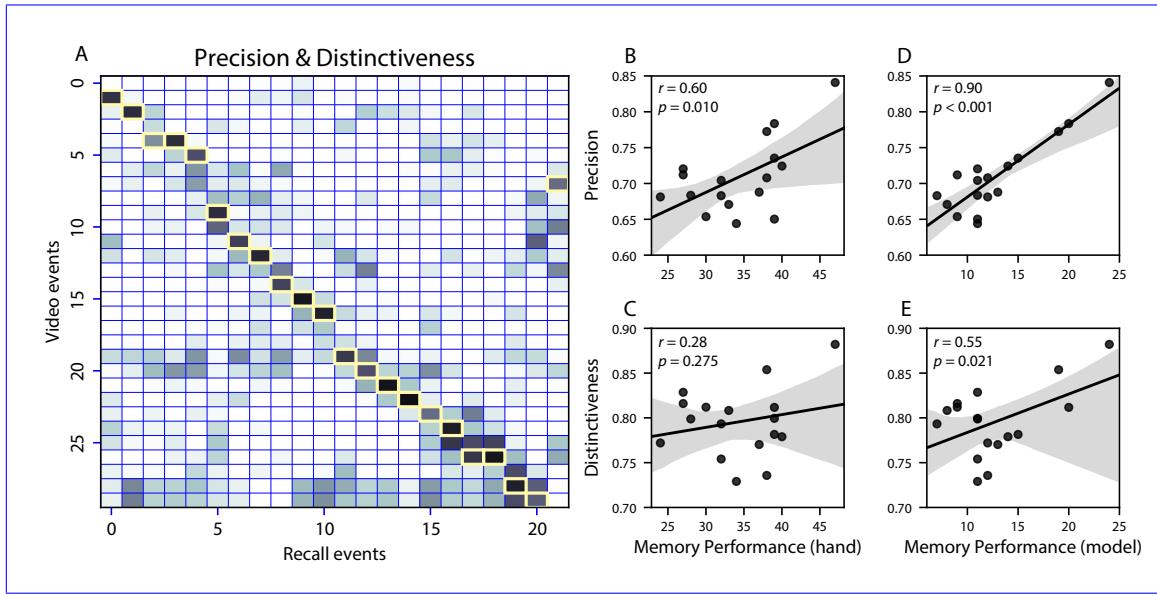


Figure 4: Novel content-based metrics of naturalistic memory: precision and distinctiveness. **A.** The video-recall correlation matrix for a representative participant (17). The yellow boxes highlight the maximum correlation in each column. The example participant's overall precision score was computed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

and the number of video events successfully remembered, as determined by our model (Pearson's $r(15) = 0.90, p < 0.001$). We also hypothesized that participants who recounted events in a more distinctive way would display better overall memory. We find that this distinctiveness score is related to our model's estimated number of recalled events (Pearson's $r(15) = 0.55, p = 0.021$), and while we do not find distinctiveness to be related to hand-annotated memory performance (Pearson's $r(15) = 0.28, p = 0.275$), this is not entirely surprising given how the hand-annotated memory scores were computed (see *Methods* and *Discussion* for details).

Further intuition for the behaviors captured by these two metrics may be gained by directly examining the content of the video and ~~participants' recalls is quantified in the corresponding~~ recalls our framework models. In Figure 5, we contrast recalls for the same video event (event 22)



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.

from two participants: one with a high precision score (P17), the other with a low precision score (P6). From the HMM-identified event boundaries, we recovered the set of annotations describing the content of an example video event (Fig. 5B), and divided them into different color-coded sections for each action or feature described. We then similarly recovered the set of sentences comprising the corresponding recall event for each of the two example participants. Because the recall sliding windows overlap heavily, and each recall event spans multiple recall timepoints (i.e., windows), we have stripped any sentences from the beginning and end that describe earlier or later video events for the sake of readability. In other words, Fig. 5C shows a subset of the full

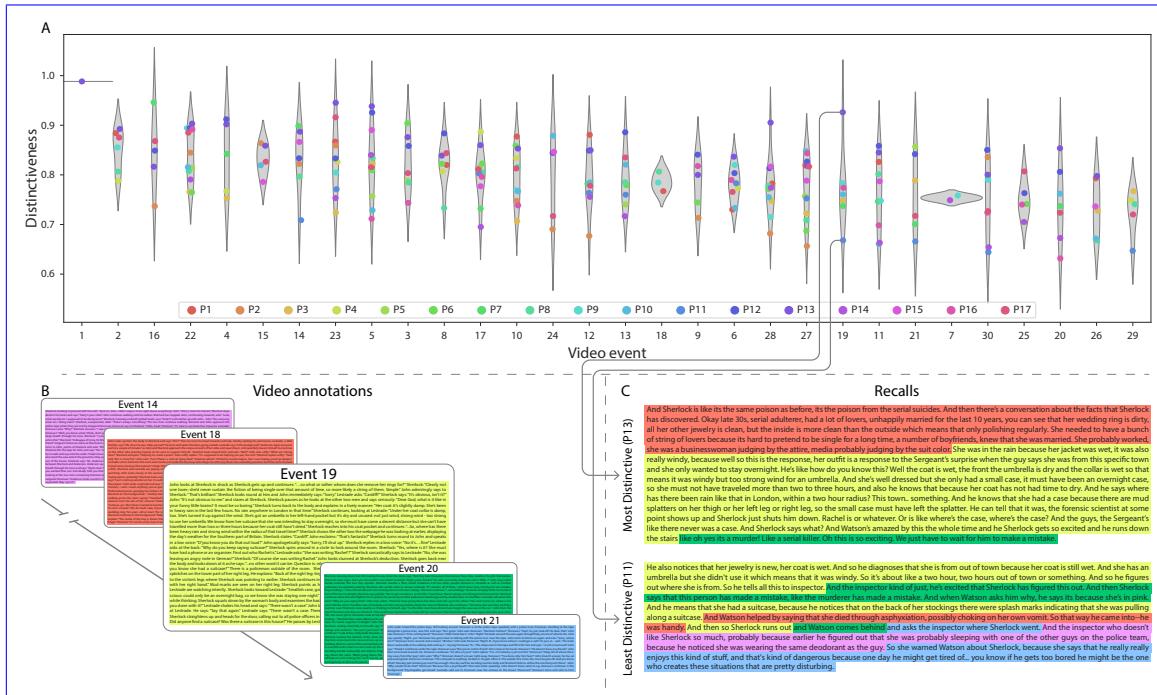


Figure 6: Distinctiveness metric reflects specificity of recall. A. Recall distinctiveness by video event. Kernel density estimates for each video event’s distribution of recall distinctiveness scores, analogous to Fig. 5A. B. The sets of “Narrative Details” video annotations (generated by Chen et al., 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.

300 recall event text, comprising sentences between the first and last descriptions of content from the
301 example video event. We then colored all words describing actions and features coded in panel
302 B by their corresponding color. Visual comparison of the transcripts reveals that the most precise
303 participant's recall both captures more of the video event's content, and does so with far more
304 detail.

305 Figure 6 similarly contrasts two example participants' recalls for a common video event (event
306 19) to illustrate the tangible differences between high and low distinctiveness scores. Here, we
307 have extracted the full set of sentences comprising the most distinctive recall event (P13) and least
308 distinctive recall event (P11) matched to the example video event (Fig. 6C). We also extracted the

309 annotations for the example video event, as well as those from each other video event whose content
310 the example participants' single recall events described (Fig. 6B). We then shaded the annotation
311 text for each video event with a different color, and shaded each word of the example participants'
312 recall text by the color of the video event it describes. The majority of the most distinctive recall
313 event text describes video event 19's content, with the first five and last one sentence describing
314 the video events immediately preceding and succeeding the current one, respectively. In contrast,
315 the least precise participant's recall for video event 19 blends the content from five separate video
316 events, does not transition between them in order, and often combines descriptions of two video
317 events' content in the same sentence.

318 The prior analyses leverage the correspondence between the 100-dimensional topic proportion
319 matrices for the video and participants' recalls to characterize recall. However, it is difficult to gain
320 deep insights into that content the content of (or relationships between) experiences and memories
321 solely by examining the topic proportion matrices these topic proportions (e.g., Figs. 2A, D) or the
322 corresponding correlation matrices (Figs. 2B, E, S4). And while we can directly examine the original
323 text underlying these topic vectors (e.g., Figs. 5, 6) to show how relationships between them reflect
324 real-world behavior, this comparison becomes prohibitively cumbersome at larger timescales.
325 To visualize the time-varying high-dimensional content in a more intuitive way (Heusser et al.,
326 2018b), we projected the topic proportions matrices onto a two-dimensional space using Uniform
327 Manifold Approximation and Projection (UMAP; ?)(UMAP; McInnes et al., 2018). In this lower-
328 dimensional space, each point represents a single video or recall event, and the distances between
329 the points reflect the distances between the events' associated topic vectors (Fig. 7). In other words,
330 events that are nearer to each other in this space are more semantically similar, and those that are
331 farther apart are less so.

332 Visual inspection of the video and recall topic trajectories reveals a striking pattern. First, the
333 topic trajectory of the video (which reflects its dynamic content; Fig. 7A) is captured nearly perfectly
334 by the averaged topic trajectories of participants' recalls (Fig. 7B). To assess the consistency of these
335 recall trajectories across participants, we asked: given that a participant's recall trajectory had
336 entered a particular location in the reduced topic space, could the position of their *next* recalled

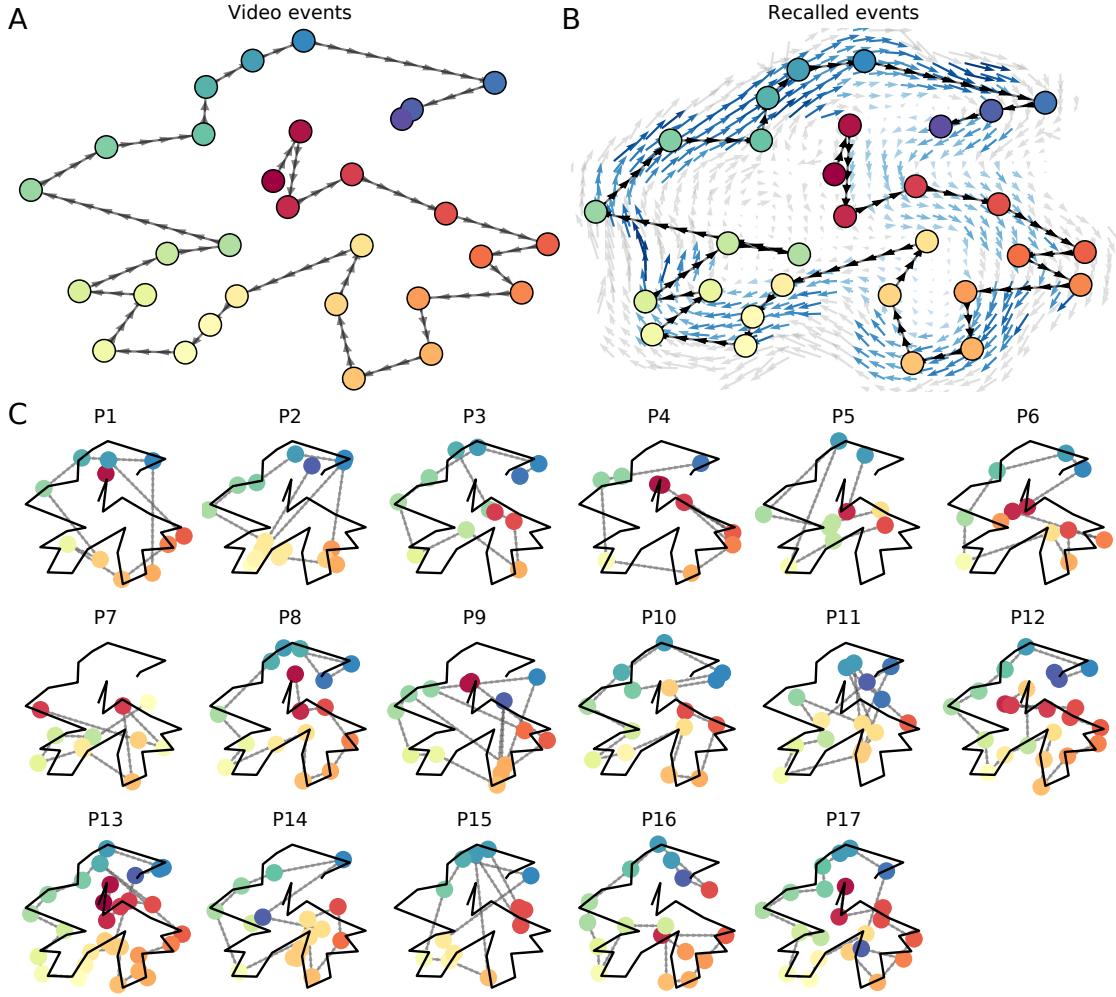


Figure 7: Trajectories through topic space capture the dynamic content of the video and recalls. All panels: the topic proportion matrices have been projected onto a shared two-dimensional space using UMAP. **A.** The two-dimensional topic trajectory taken by the episode of *Sherlock*. Each dot indicates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

337 event be predicted reliably? For each location in the the reduced topic space, we computed the set of
338 line segments connecting successively recalled events (across all participants) that intersected that
339 location (see *Methods* for additional details). We then computed (for each location) the distribution
340 of angles formed by the lines defined by those line segments and a fixed reference line (the x -
341 axis). Rayleigh tests revealed the set of locations in topic space at which these across-participant
342 distributions exhibited reliable peaks (blue arrows in Fig. 7B reflect significant peaks at $p < 0.05$,
343 corrected). We observed that the locations traversed by nearly the entire video trajectory exhibited
344 such peaks. In other words, participants exhibited similar trajectories that also matched the
345 trajectory of the original video (Fig. 7C). This is especially notable when considering the fact that
346 the number of events participants recalled (dots in Fig. 7C) varied considerably across people, and
347 that every participant used different words to describe what they had remembered happening in
348 the video. Differences in the numbers of remembered events appear in participants' trajectories
349 as differences in the sampling resolution along the trajectory. We note that this framework also
350 provides a means of detangling disentangling classic "proportion recalled" measures (i.e., the
351 proportion of video events referenced described in participants' recalls) from participants' abilities
352 to recapitulate the full-gist overall unfolding of the original video's content (i.e., the similarity in
353 the shape between the shapes of the original video trajectory and that defined by each participant's
354 recounting of the video).

355 In addition to the more "holistic" measure of memory described in the previous section, our
356 framework also affords the ability to drill down to individual words and quantify how each word
357 relates to the memorability of each event. The results displayed in Figures 3C and 5A suggest that
358 certain events were remembered better than others. Given this, we next asked asked whether the
359 events were generally remembered well or poorly tended to reflect particular content. Because our
360 analysis framework projects the dynamic video content and participants' recalls onto a shared topic
361 into a shared space, and because the dimensions of that space are known (i.e., each topic dimension
362 is a set represent topics (which are, in turn, sets of weights over words in the vocabulary; Fig. S2), we
363 can examine the topic trajectories to understand which specific content was remembered well (or
364 poorly). For each video event, we can ask: what was the average correlation (across participants)

365 between the video event's topic vector and the closest matching recall event topic vectors from
366 each participant? This yields a single correlation coefficient for each video event, describing how
367 closely participants' recalls of the event tended to reliably capture its content are able to recover
368 the weighted combination of words that make up any point (i.e., topic vector) in this space. We
369 first computed the average precision with which participants recalled each of the 30 video events
370 (Fig. 8A). (We also examined how different comparisons between each video event's topic vector
371 and the corresponding recall event topic vectors related to hand-annotated characterizations of
372 memory performance; see *Supporting Information*). Given this summary of which events were
373 recalled reliably (or not), we next asked whether the better-remembered or worse-remembered
374 events tended to reflect particular topics. We note that this result is analogous to a serial position
375 curve created from our continuous recall quality metric). We then computed a weighted average
376 of the topic vectors for each video event, where the weights reflected how reliably each event
377 was recalled. To visualize the result, we created a "wordle" image (Mueller et al., 2018) where
378 words weighted more heavily by better-remembered topics appear in a larger font (Fig. 8B, green
379 box). Events Across the full video, content that reflected topics weighting heavily on characters
380 like "Sherlock" and "John" (i.e., the main characters) and locations like "221b Baker Street" (i.e.,
381 a major recurring location necessary to convey the central focus of the video (e.g., the names of
382 the two main characters, "Sherlock" and "John", and the address of the flat that Sherlock and
383 John share a major recurring location, "221B Baker Street") were best remembered. An analogous
384 analysis revealed which themes were poorly remembered. Here in computing the weighted
385 average over events' topic vectors, we weighted each event in *inverse* proportion to how well it
386 was remembered (Fig. 8B, red box). This revealed that events with The least well-remembered
387 video content reflected information not necessary to later convey a general summary of the video,
388 such as the proper names of relatively minor characters such as (e.g., "Mike," "Jeffrey," and "Molly,"
389 as well as less integral plot "Molly," and "Lestrade") and locations (e.g., "hospital" and "office")
390 were least well-remembered. This suggests that what is retained in memory are the major plot
391 elements (i.e., the overall "gist" of what happened), whereas the more minor details are prone to
392 pruning. St. Bartholomew's Hospital").

393 In addition to constructing overall summaries, assessing the video and recall topic vectors from
394 individual recalls can provide further insights. Specifically, for any given event we can construct
395 A similar result emerged from assessing the topic vectors for individual video and recall
396 events (Fig. 8C). Here, for each of the three best- and worst-remembered video events, we have
397 constructed two wordles: one from the original video event's topic vector τ (left) and a second from
398 the average topic vectors produced by all participants' recalls of that event. We can then examine
399 those wordles visually to gain an intuition for which aspects of the video event were recapitulated
400 in participants' recalls of that event. Several example wordles are displayed in Figure 8C (wordles
401 from the recall topic vector for that event (right). The three best-remembered events are (circled
402 in green; wordles from the) correspond to scenes important to the central plot-line: a mysterious
403 figure spying on John in a phone booth; John meeting Sherlock at Baker St. to discuss the murders;
404 and Sherlock laying a trap to catch the killer. Meanwhile, the three worst-remembered events are
405 (circled in red) . Using wordles to visually compare the topical content of each video event and
406 the (average) corresponding recall event reveals the specific content from the specific events that
407 is reliably retained in the transformation into memory (green events) or not (red events) reflect
408 scenes that are non-essential to summarizing the narrative's structure: the video of singing cartoon
409 characters participants viewed prior to the main episode; John asking Molly about Sherlock's habit
410 of over-analyzing people; and Sherlock noticing evidence of Anderson's and Donovan's affair.

411 **Transforming experience into memory.** A. Average correlations (across participants) between
412 the topic vectors from each video event and the closest matching recall events. Error bars
413 denote bootstrap-derived across-participant 95% confidence intervals. The stars denote the three
414 best-remembered events (green) and worst-remembered events (red). B. Wordles comprising the
415 top 200 highest-weighted words reflected in the weighted-average topic vector across video events.
416 Green: video events were weighted by how well the topic vectors derived from recalls of those
417 events matched the video events' topic vectors (Panel A). Red: video events were weighted by the
418 inverse of how well their topic vectors matched the recalled topic vectors. C. The set of all video and
419 recall events is projected onto the two-dimensional space derived in Figure 7. The dots outlined
420 in black denote video events (dot size reflects the average correlation between the video event's

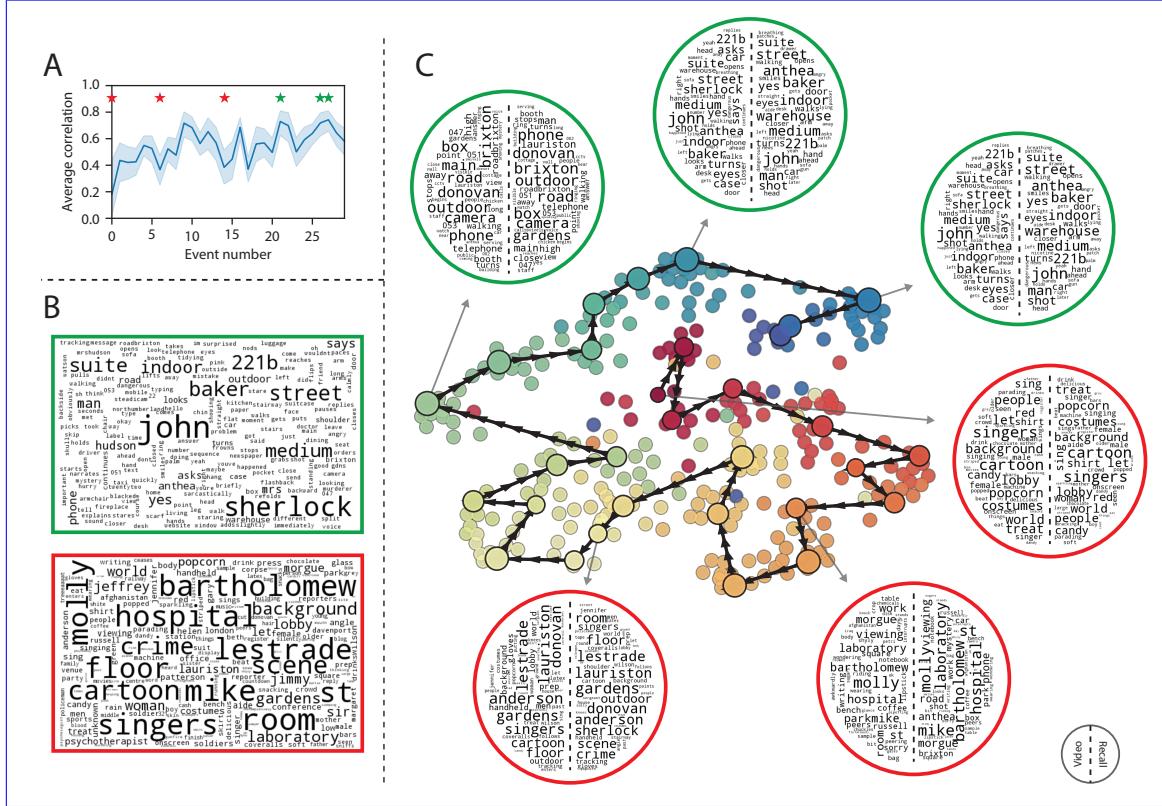


Figure 8: Language used in the most and least memorable events. **A.** Average precision (video event-recall event topic vector correlation) across participants for each video event. Error bars denote bootstrap-derived across-participant 95% confidence intervals. The stars denote the three best-remembered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

421 topic vector and the topic vectors from the closest matching recalled events from each participant;
422 bigger dots denote stronger correlations). The dots without black outlines denote recalled events.
423 All dots are colored using the same scheme as Figure 7A. Wordles for several example events are
424 displayed (green: three best remembered events; red: three worst remembered events). Within
425 each circular wordle, the left side displays words associated with the topic vector for the video
426 event, and the right side displays words associated with the (average) recall event topic vector,
427 across all recall events matched to the given video event.

428 The results thus far inform us about which aspects of the dynamic content in the episode
429 participants watched were preserved or altered in participants' memories of the episode. We next
430 carried out a series of analyses aimed at understanding which brain structures might implement
431 these processes. In one analysis facilitate these preservations and transformations between the
432 external world and memory. In the first analysis, we sought to identify which brain structures
433 brain structures that were sensitive to the dynamic unfolding of the video's dynamic content, as
434 characterized by its topic trajectory. Specifically, we used a searchlight procedure to identify the
435 extent to which each cluster of voxels exhibited a timecourse (as the clusters of voxels whose activity
436 patterns displayed a proximal temporal correlation structure (as participants watched the video)
437 whose temporal correlation matrix matched the temporal correlation matrix matching that of the
438 original video's topic proportion matrix proportions (Fig. 2B). As shown in Figure 9A, the analysis
439 revealed a network of regions including bilateral frontal cortex and cingulate cortex, suggesting
440 that these regions may play a role in maintaining information relevant to the narrative structure of
441 the video (see Methods for additional details). In a second analysis, we sought to identify which
442 brain structures' responses (while viewing the video brain structures whose responses (during
443 video viewing) reflected how each participant would later recall structure their recounting of the
444 video. We used an analogous searchlight procedure to identify clusters of voxels whose proximal
445 temporal correlation matrices reflected the temporal correlation matrix matched that of the topic
446 proportions for each individual's recalls recall (Figs. 2D, S4, 9B; see Methods for additional details).
447 To ensure our searchlight procedure identified regions specifically sensitive to the temporal structure
448 of the video or recalls (i.e., rather than those with a temporal autocorrelation length similar to that

449 of the video/recalls), we performed a phase shift-based permutation correction (see *Methods* for
450 additional details). As shown in Figure 9B, the C, the video-driven searchlight analysis revealed
451 a distributed network of regions including the ventromedial prefrontal cortex (vmPFC), anterior
452 eingulate cortex, and right medial temporal lobe (rMTL), suggesting that these regions may play
453 a role in transforming each individual that may play a role in processing information relevant to
454 the narrative structure of the video. Similarly, the recall-driven searchlight analysis revealed a
455 second network of regions (Fig. 9D) that may facilitate a person-specific transformation of one's
456 experience into memory. In identifying regions whose responses to ongoing experiences reflect
457 how those experiences will be remembered later, this latter analysis extends classic *subsequent*
458 *memory analyses* (e.g., Paller and Wagner, 2002) to domain of naturalistic stimuli.

459 The searchlight analyses described above yielded two distributed networks of brain regions,
460 whose activity timecourses mirrored to the temporal structure of the video (Fig. 9C) or participants'
461 eventual recalls (Fig. 9D). We next sought to gain greater insight into the structures and functional
462 networks our results reflected. To accomplish this, we performed an additional, exploratory
463 analysis using Neurosynth (Yarkoni et al., 2011). Given an arbitrary statistical map as input,
464 Neurosynth performs a massive automated meta-analysis, returning a ranked list of terms reported
465 in papers with similar significance maps. We ran Neurosynth on the significance maps for the video-
466 and recall-driven searchlight analyses. These maps, along with the 10 terms with maximally similar
467 meta-analysis images identified by Neurosynth are shown in Figure 10.

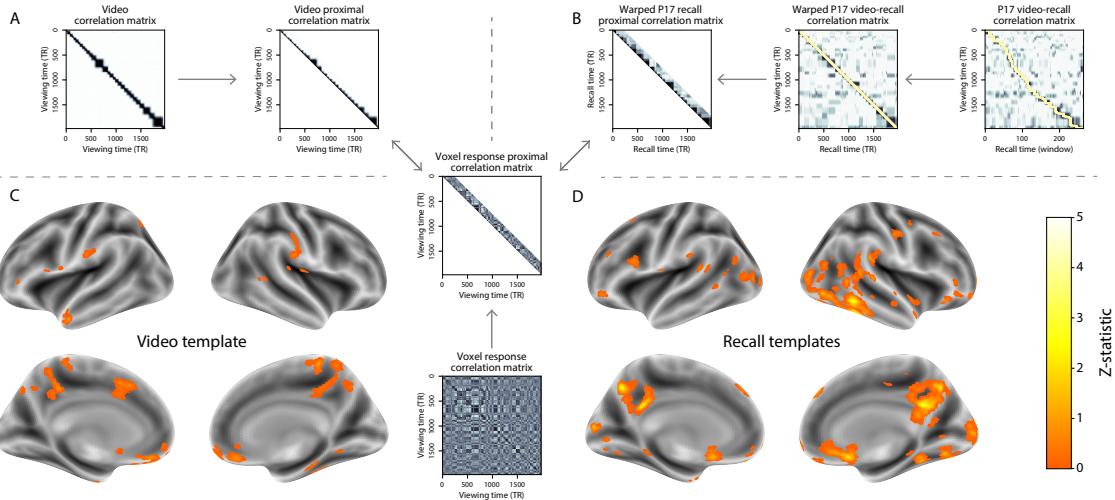


Figure 9: Brain structures that underlie the transformation of experience into memory. A. We searched for regions whose responses (as participants watched isolated the video) matched proximal diagonals from the temporal correlation matrix upper triangle of the video topic proportions. These regions are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

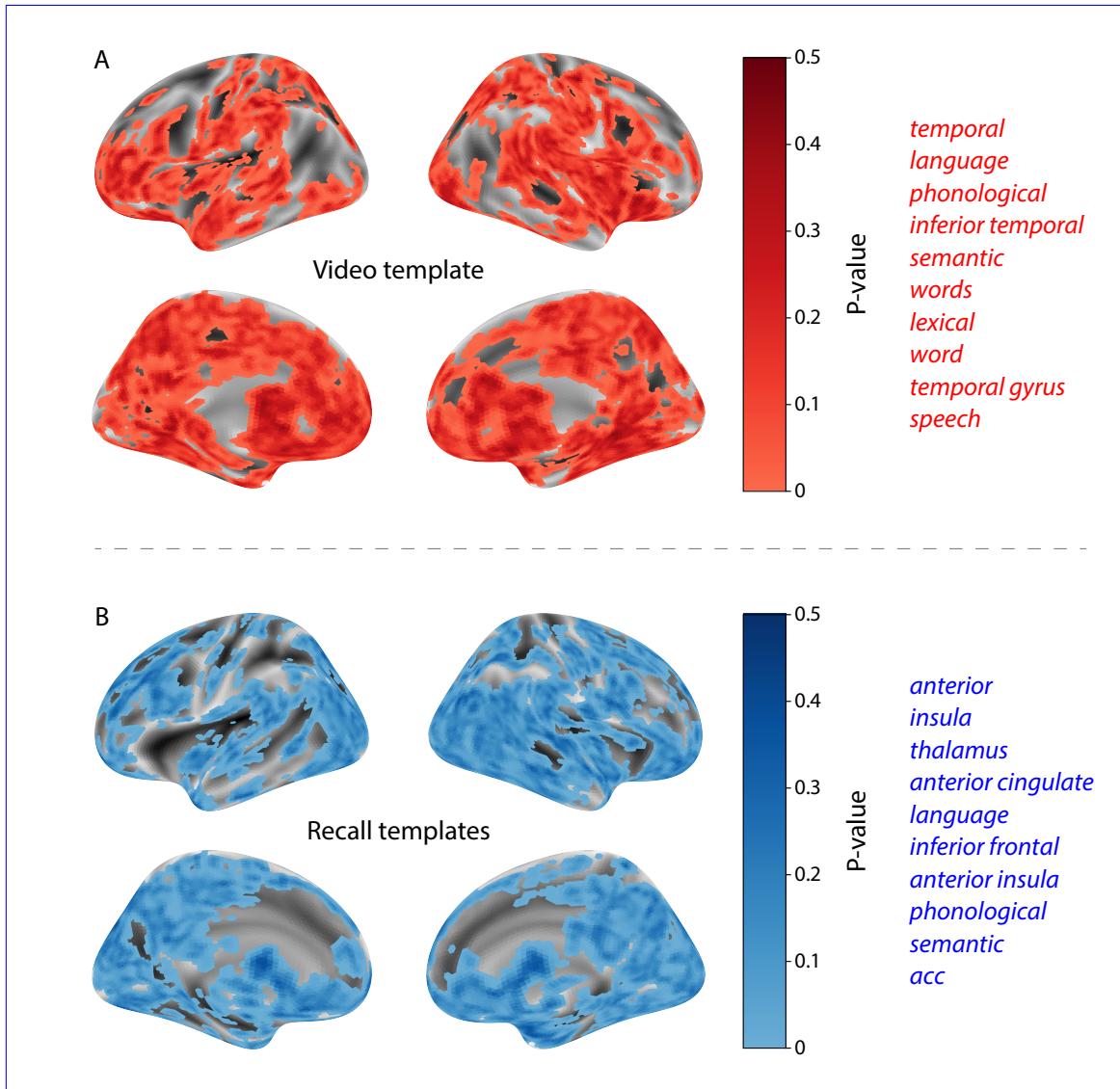


Figure 10: Decoding distributed statistical maps via Neurosynth meta-analyses. **A.** Video-searchlight significance and top 10 decoded terms. We constructed a map of the permutation-derived p -values for the video-driven searchlight analysis (Fig. 9A, C) at each voxel with a positive permutation-derived 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.

468 **Discussion**

469 Our work casts remembering as reproducing (behaviorally and neurally) the topic trajectory,
470 or “gist, ” of the original shape, of an experience. This view draws inspiration from prior
471 work aimed at elucidating the neural and behavioral underpinnings of how we process dy-
472 namic naturalistic experiences and remember them later. One approach to identifying neural
473 responses to naturalistic stimuli (including experiences) entails building a model of the stimu-
474 lus and searching for brain regions whose responses are consistent with the model. In prior
475 work, a series of studies from Uri Hasson’s group (Lerner et al., 2011; Simony et al., 2016; Chen
476 et al., 2017; Baldassano et al., 2017; Zadbood et al., 2017) have extended this approach with a
477 clever twist.—Rather: rather than building an explicit stimulus model, these studies instead
478 search for brain responses (while experiencing the stimulus) that are reliably similar across in-
479 dividuals. So called *inter-subject correlation* (ISC) and *inter-subject functional connectivity* (ISFC)
480 analyses effectively treat other people’s brain responses to the stimulus as a “model” of how its
481 features change over time. By contrast, in our present workwe used topic models and HMMs
482 , we use topic models to construct an explicit stimulus model content model directly from the
483 stimulus (i.e., the topic trajectory of the video). When we searched for brain structures whose
484 responses are consistent with the video’s topic trajectory, we identified a network of structures that
485 overlapped strongly with the “long temporal receptive window” network reported by the Hasson
486 group (e.g., compare our Fig. 9A with the map of long temporal receptive window voxels in Lerner et al., 2011)
487 . This provides support for the notion that part of the long temporal receptive window network may
488 be maintaining an explicit model Projecting each participant’s recall into a space shared by both the
489 stimulus and other participants then allows us to compare recalls both directly to the stimulus and
490 to each other. Similarly, prior work introducing the use of HMMs to discover latent event structure
491 in naturalistic stimuli and recall (Baldassano et al., 2017) used between-subjects cross-validation
492 to identify event boundaries shared across participants, and between stimulus and recall. Our
493 framework allows us to break from the restriction of a common, shared event-timeseries and
494 identify the unique resolution of each participant’s recall event structure, and how that may differ

495 from the video and that of other participants.

496 While a large number of language models exist (e.g., WAS, LSA, word2vec, universal sentence encoder; Steyvers et al.
497 , here we use latent dirichlet allocation (LDA)-based topic models for a few reasons. First, topic
498 models capture the *essence* of a text passage devoid of the specific set and order of words used. This
499 was an important feature of our model since different people may accurately recall a scene using
500 very different language. Second, words can mean different things in different contexts (e.g. "bat"
501 may be the act of hitting a baseball, the object used for that action, or as a flying mammal). Topic
502 models are robust to this, allowing words to exist as part of multiple topics. Last, topic models
503 provide a straightforward means to recover the weights for the particular words comprising a topic,
504 enabling easy interpretation of an event's contents (e.g. Fig. 8). Other models such as Google's
505 Universal Sentence Encoder offer a context-sensitive encoding of text passages, but the encoding
506 space is complex and non-linear, and thus recovering the original words used to fit the model is
507 not straightforward. However, it's worth pointing out that our framework is divorced from the
508 particular choice of language model. Moreover, many of the aspects of our framework could be
509 swapped out for other choices. For example, the language model, the timeseries segmentation
510 model and the video-recall matching function could all be customized for the particular problem.
511 Indeed for some problems, recovery of the particular recall words may not be necessary, and thus
512 other text-modeling approaches (such as universal sentence encoder) may be preferable. Future
513 work will explore the influence of particular model choices on the framework's efficacy.

514 In extending classical free recall analyses to our naturalistic memory framework, we recovered
515 two patterns of recall dynamics central to list-learning studies: a heightened probability of initiating
516 recall with the first presented "item" (in our case, video events; Fig. 3A) and a strong bias toward
517 transitioning from recalling a given event to recalling the one immediately following it (Fig. 3B).
518 However, equally noteworthy are the typical free recall results *not* recovered in these analyses,
519 as each highlights a fundamental difference between the list-learning paradigm and naturalistic
520 memory paradigms like the one employed in the present study. The most noticeable departure
521 from hallmark free recall dynamics in these findings is the apparent lack of a serial position effect in
522 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). 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, 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 reflected in our findings’ second departure from typical free recall dynamics: a lack of increased probability of first recall for end-of-sequence events (Fig. 3A).

Because they disregard presentation order-dependent variability in the stimulus content, analyses such as those in Figure 3 enable a more sensitive analysis of presentation order-dependent temporal dynamics in free recall. Yet by the same token, they paint a wholly incomplete picture of memory for naturalistic episodes. In an attempt to address this shortcoming, we have developed a framework in the present study that characterizes the explicit semantic content of the stimulus dynamics. When we performed a similar analysis after swapping out the video’s topic trajectory

551 with the recall topic trajectories of each individual participant, this allowed us to identify and
552 subsequent recalls. However, sensitivity to stimulus and recall content introduces a new challenge:
553 distinguishing between levels of recall quality for a stimulus (e.g., an event) that is considered to
554 have been “remembered.” When modeling memory in an experimental setting, recall quality for
555 individual events is often cast as binary (e.g., a given list item was simply either remembered or
556 not remembered). Various models of memory (e.g., Yonelinas, 2002) attempt to improve upon this
557 by including confidence ratings, rendering this binary judgement instead categorical. To better
558 evaluate naturalistic memory quality, we introduce a continuous metric (*precision*), which reflects
559 the level of completeness of a participant’s recall for a feature-rich experience. Additionally, recall
560 quality for a single event is typically assessed independently from that for all other events (e.g., it
561 is difficult to “compare” a participant’s binary recall success for list item 1 to that of list item 10).
562 The second novel metric we introduce (*distinctiveness*) is based on analyzing of the correlational
563 structure of an individual’s full set of recall events, and reflects the specificity of their memory
564 for a single experienced event. We find that both of these metrics relate to the overall number of
565 video events participants successfully recalled, and that our precision metric additionally relates to
566 Chen et al. (2017)’s hand-annotated memory memory scores. Though we do not find participants’
567 average recall distinctiveness related to the hand-annotated memory scores, this is not entirely
568 surprising given the divergence of behavior they capture. In hand-scoring each participant’s
569 verbal recall for each of 50 (manually-delimited) scenes, “[a] scene was counted as recalled if the
570 participant described any part of the scene” (Chen et al., 2017). In other words, both an extensive
571 description of a scene’s content and a brief mention of some subset of its content were (binarily)
572 considered equally successful recalls. By contrast, we identify the event structure in participants’
573 recalls in an unsupervised manner, independent of the video event-timeseries, prior to mapping
574 between video and recall content. Our HMM-based event-segmentation produces boundaries
575 between timepoints where the topic proportions shift in a substantial way, and because a small
576 handful of words is unlikely to contribute significantly to the topic proportions for any sliding
577 window, such brief scene descriptions will most often not begat a sufficiently large shift in the
578 resulting topic proportions for the HMM to identify an event boundary. Instead, they will be

579 grouped with a neighboring event, consequently lowering that event's distinctiveness score and
580 by extension, the participant's overall distinctiveness score. This is in essence the qualitative
581 difference between distinctive and indistinctive recall, and reflects the comparison shown in Figure
582 6C. Intriguingly, prior studies show that pattern separation, or the ability to cleanly discriminate
583 between similar experiences, is impaired in many cognitive disorders as well as natural aging
584 (Stark et al., 2010; Yassa et al., 2011; Yassa and Stark, 2011). Future work might explore whether
585 and how these metrics compare between cognitively impoverished groups and healthy controls.

586 In the analyses outlined in Figure 9, we identified two networks of brain regions whose re-
587 sponds (as the participants viewed the video) reflected how during video viewing were consistent
588 with the temporal structure of the video and recall topic trajectories, respectively. The network
589 identified by the video trajectory would be transformed in memory (as reflected by the recall topic
590 trajectories). The analysis revealed that the rMTL and vmPFC analysis included the ventromedial
591 prefrontal cortex, left anterior temporal lobe, superior parietal and dorsal anterior cingulate cortex.
592 The network from the video-recall trajectory analysis also included the ventromedial prefrontal and
593 superior parietal cortices, in addition to the posterior medial cortex (PMC) and the inferior temporal
594 regions. Notably, Chen et al. (2017) also observed the PMC in a number of analyses including one
595 that searched for regions whose activity patterns during encoding were reinstated during free
596 recall. The PMC has been consistently identified in studies involving stimuli with meaningfully
597 structured events ?. Further, the PMC is part of the "posterior medial" system, a network of brain
598 regions thought to represent situation models Zacks et al. (2007) in support of memory, spatial
599 navigation and social cognition (Ranganath and Ritchey, 2012). Given that we constructed our
600 video-recall searchlight model to capture temporal structure in the episode's semantic content
601 (and how one's later recall aligns with that structure), we speculate that the PMC may play a role
602 in this person-specific transformation from experience into memory. The role of the MTL in episodic
603 memory encoding has been well-reported (e.g., Paller and Wagner, 2002; Davachi et al., 2003; ?, Davachi, 2006)
604 . Prior work has also implicated the medial prefrontal cortex in representing "schema" knowledge (i.e., general knowledge
605 . Integrating across our study and this prior work, one interpretation is that the person-specific
606 transformations mediated (or represented) by the rMTL and vmPFC may reflect schema knowledge

607 being leveraged, formed, or updated, incorporating ongoing experience into previously acquired
608 knowledge, constructing mnemonic events from meaningfully structured experiences.

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

618

619 Our second searchlight analysis identified a largely separate network of regions (Fig. 9D)
620 whose patterns of activity as participants viewed the video reflected the idiosyncratic structure
621 of each individual's later recall. Decoding the associated significance map yielded a set of terms
622 that primarily reflected names of specific structural regions (such as "thalamus," "anterior insula,"
623 "anterior cingulate" and "inferior frontal"; Fig. 10B). Interestingly, these regions share membership
624 in a common, large-scale functional network (termed the "salience network") involved in detecting
625 and processing affective cues. In particular, the latter three regions have been implicated in
626 functions relevant to assigning personal meaning to an experience, including: ascribing subjective
627 value to raw, sensory input (?); modulating semantic and phonological processing in response
628 to personally salient stimuli (?); and directing and reallocating attention and working memory
629 resources towards the most relevant stimuli (Menon and Uddin, 2010). This suggests that the
630 network of structures displayed in Figure 9D may play a role in transforming and restructuring
631 ongoing experiences through the lens of one's own personal values as they are encoded in memory.

632

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

635 classroom evaluation tools involve simply computing the proportion of correctly answered exam
636 questions. Our work indicates that this approach is only loosely related to what educators might
637 really want to measure: how well did the students understand the key ideas presented in the
638 course? ~~One could apply the computational framework we developed to construct topic trajectories~~
639 ~~for the video and participants' recalls~~ Under this typical framework of assessment, the same exam
640 score of 50% could be ascribed to two very different students: one who attended the full course
641 but struggled to learn more than a broad overview of the material, and one who attended only half
642 of the course but understood the material perfectly. Instead, one could apply our computational
643 framework to build explicit content models of the course material and exam questions. This
644 approach would provide a more nuanced and specific view into which aspects of the material
645 students had learned well (or poorly). In clinical settings, memory measures that incorporate such
646 explicit content models might also provide more direct evaluations of patients' memories.

647 Methods

648 Experimental design and data collection

649 Data were collected by Chen et al. (2017). In brief, participants (~~n = 17~~n = 22) viewed the first 48
650 minutes of "A Study in Pink", the first episode of the BBC television series *Sherlock*, while fMRI
651 volumes were collected (TR = 1500 ms). Participants were pre-screened to ensure they had never
652 seen any episode of the show before. The stimulus was divided into a 23 min (946 TR) and a 25 min
653 (1030 TR) segment to mitigate technical issues related to the scanner. After finishing the clip,
654 participants were instructed to (quoting from Chen et al., 2017) "describe what they recalled of the
655 [episode] in as much detail as they could, to try to recount events in the original order they were
656 viewed in, and to speak for at least 10 minutes if possible but that longer was better. They were told
657 that completeness and detail were more important than temporal order, and that if at any point
658 they realized they had missed something, to return to it. Participants were then allowed to speak
659 for as long as they wished, and verbally indicated when they were finished (e.g., 'I'm done')."Five

660 participants were dropped from the original dataset due to excessive head motion (2 participants),
661 insufficient recall length (2 participants), or falling asleep during stimulus viewing (1 participant),
662 resulting in a final sample size of $n = 17$. For additional details about the experimental procedure
663 and scanning parameters, see Chen et al. (2017). The experimental protocol was approved by
664 Princeton University's Institutional Review Board.

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

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

678 Data and code availability

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

681 Statistics

682 All statistical tests we performed performed in the behavioral analyses were two-sided. All
683 statistical tests performed in the neural data analyses were two-sided, except for the permutation-based
684 thresholding, which was one-sided. In this case, we were specifically interested in identifying

685 voxels whose activation time series reflected the temporal structure of the video and recall
686 trajectories to a *greater* extent than that of the phase-shifted trajectories.

687 Modeling the dynamic content of the video and recall transcripts

688 Topic modeling

689 The input to the topic model we trained to characterize the dynamic content of the video com-
690 prised 998 hand-generated annotations of each of 1000 short (mean: 2.96s) scenes spanning the
691 video clip (generated by Chen et al., 2017). The features included: narrative details (a sentence or
692 two describing what happened in that scene); whether the scene took place indoors or outdoors;
693 names of any characters that appeared in the scene; name(s) of characters in camera focus; name(s)
694 of characters who were speaking in the scene; the location (in the story) that the scene took place;
695 camera angle (close up, medium, long, top, tracking, over the shoulder, etc.); whether music was
696 playing in the scene or not; and a transcription of any on-screen text (Chen et al., 2017 generated
697 1000 annotations total; we removed two referring to the break between the first and second scan
698 sessions, during which no fMRI data was collected). We concatenated the text for all of these
699 the annotated features within each segment, creating a “bag of words” describing each scene
700 —and performed some minor preprocessing (e.g., stemming possessive nouns and removing
701 punctuation). We then re-organized the text descriptions into overlapping sliding windows span-
702 ning (up to) 50 scenes each. In other words, the first text sample comprised the combined text
703 from the first 50 scenes we created a “context” for each scene comprising the text descriptions of
704 the preceding 25 scenes, the present scene, and the following 24 scenes. To model the “context”
705 for scenes near the beginning and end of the video (i.e., 1–50), the second comprised the text from
706 scenes 2–51, and so on. within 25 scenes of the beginning or end), we created overlapping sliding
707 windows that grew in size from one scene to the full length, then similarly tapered their length at
708 the end. This additionally ensured that each scene’s content was represented in the text corpus an
709 equal number of times.

710 We trained our model using these overlapping text samples with scikit-learn (version 0.19.1;

711 Pedregosa et al., 2011), called from our high-dimensional visualization and text analysis software,
712 HyperTools (Heusser et al., 2018b). Specifically, we ~~use~~used the CountVectorizer class to trans-
713 form the text from each ~~scene~~window into a vector of word counts (using the union of all words
714 across all scenes as the “vocabulary,” excluding English stop words); this ~~yields a number-of-scenes~~yielded a number-of-windows by number-of-words *word count* matrix. We then ~~use~~used the
715 LatentDirichletAllocation class (topics=100, method='batch') to fit a topic model (Blei et al.,
716 2003) to the word count matrix, yielding a ~~number-of-scenes~~(1000)~~number-of-windows~~(1047) by
717 number-of-topics (100) *topic proportions* matrix. The topic proportions matrix describes ~~which~~the
718 ~~gradually evolving~~ mix of topics (latent themes) ~~is~~ present in each scene. Next, we transformed the
719 topic proportions matrix to match the 1976 fMRI volume acquisition times. ~~For each fMRI volume,~~
720 ~~we took the topic proportions from whatever scene was displayed for most of that volume's 1500 ms~~
721 ~~acquisition time. This yielded a new~~We assigned each topic vector to the timepoint (in seconds)
722 ~~midway between the beginning of the first scene and the end of the last scene in its corresponding~~
723 ~~sliding text window. By doing so, we warped the linear temporal distance between consecutive~~
724 ~~topic vectors to align with the inconsistent temporal distance between consecutive annotations~~
725 ~~(whose durations varied greatly). We then rescaled these timepoints to 1.5s TR units, and used~~
726 ~~linear interpolation to estimate a topic vector for each TR. This resulted in a~~number-of-TRs (1976)
727 by number-of-topics (100) ~~topic proportions~~ matrix.

728
729 We created similar topic proportions matrices using hand-annotated transcripts of each par-
730 ticipant’s recall of the video (annotated by Chen et al., 2017). We tokenized the transcript into a
731 list of sentences, and then re-organized the list into overlapping sliding windows spanning (up
732 ~~to~~) 10 sentences each; ~~in turn~~, analogously to how we parsed the video annotations. In turn, we
733 transformed each window’s sentences into a word count vector (using the same vocabulary as for
734 the video model). ~~We~~, then used the topic model already trained on the video scenes to compute
735 the most probable topic proportions for each sliding window. This yielded a ~~number-of-sentences~~number-of-windows (range: ~~68–294~~83–312) by number-of-topics (100) topic proportions matrix
736 ~~for~~ for each participant. These reflected the dynamic content of each participant’s recalls. Note:
737 for details on how we selected the video and recall window lengths and number of topics, see

739 Supporting Information and Figure S1.

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

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

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

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

750 **Prior studies on narrative structure and processing have shown that we both perceive and internally**
751 **represent the world around us at multiple, hierarchical timescales (e.g., Hasson et al., 2008; Lerner et al., 2011; Hasson et**
752 **). However, for the purposes of our framework, we sought to identify the single timeseries of**
753 **event-representations that is emphasized most heavily in the temporal structure of the video and**
754 **of each participant's recall. We quantified this as the set of K states that maximized the similarity**
755 **between topic vectors for timepoints comprising each state, while minimizing the similarity**
756 **between topic vectors for timepoints across different states. Specifically, we computed (for each**
757 **matrix)**

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

758 where a was the **average correlation between the topic vectors of timepoints within the same state;**
759 **distribution of within-state topic vector correlations, and** b was the **average correlation between**
760 **the topic vectors of timepoints within different states; and** α was a regularization parameter that we

761 set to 5 times the window length (i.e., 250 scenes for the video topic trajectory and distribution of
762 across-state topic vector correlations). We computed the first Wasserstein distance (W_1 ; also known
763 as “earth mover’s distance”; Dobrushin, 1970; Ramdas et al., 2017) between these distributions for
764 a large range of possible K -values (range [2, 50 sentences for the recall topic trajectories]), and
765 selected the K that yielded the maximum value. Figure 2B displays the event boundaries returned
766 for the video, and Figure S4 displays the event boundaries returned for each participant’s recalls.
767 See Figure S6 for the optimization functions for the video and recalls. After obtaining these event
768 boundaries, we created stable estimates of each topic proportions matrix—the content represented
769 in each event by averaging the topic vectors within each event across timepoints between each pair
770 of event boundaries. This yielded a number-of-events by number-of-topics matrix for the video
771 and recalls from each participant.

772 Naturalistic extensions of classic list-learning analyses

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

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

788 **Probability of first recall (PFR).** PFR curves (Welch and Burnett, 1924; Postman and Phillips, 1965; Atkinson and Shiffman, 1973) reflect the probability that an item will be recalled first as a function of its serial position during encoding. To carry out this analysis, we initialized a number-of-participants (17) by number-of-video-events (30) matrix of zeros. Then for each participant, we found the index of the video event that was recalled first (i.e., the video event whose topic vector was most strongly correlated with that of the first recall event) and filled in that index in the matrix with a 1. Finally, we averaged over the rows of the matrix, resulting in a 1 by 30 array representing the proportion of participants that recalled an event first, as a function of the order of the event's appearance in the video (Fig. 3A).

796 **Lag conditional probability curve (lag-CRP).** The lag-CRP curve (Kahana, 1996) reflects the probability of recalling a given item after the just-recalled item, as a function of their relative encoding positions (or *lag*). In other words, a lag of 1 indicates that a recalled item was presented immediately after the previously recalled item, and a lag of -3 indicates that a recalled item came 3 items before the previously recalled item. For each recall transition (following the first recall), we computed the lag between the current recall event and the next recall event, normalizing by the total number of possible transitions. This yielded a number-of-participants (17) by number-of-lags (-29 to +29; 61 lags total) matrix. We averaged over the rows of this matrix to obtain a group-averaged lag-CRP curve (Fig. 3B).

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

814 order appearance in the video (Fig. 3C).

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

824 **Semantic clustering scores.** Semantic clustering describes a participant's tendency to recall
825 semantically similar presented items together in their recall sequences. Here, we used the topic
826 vectors for each event as a proxy for its semantic content. Thus, the similarity between the semantic
827 content for two events can be computed by correlating their respective topic vectors. For each recall
828 event transition, we sorted all not-yet-recalled events according to how correlated the topic vector
829 of the closest-matching video event was to the topic vector of the closest-matching video event to the
830 just-recalled event. We then computed the percentile rank of the observed next recall. We averaged
831 these percentile ranks across all of the participant's recalls to obtain a single semantic clustering
832 score for the participant.

833 **Novel naturalistic memory metrics**

834 **Precision.** We tested whether participants who recalled more events were also more *precise* in
835 their recollections. For each participant, we computed the average correlation between the topic
836 vectors for each recall event and those of its closest-matching video event. This gave a single value
837 per participant representing the average precision across all recalled events. We then correlated
838 these values with both hand-annotated and model-derived (i.e., the number of unique video events

839 matched by a participant's recall events) memory performance.

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

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

846 where $\bar{c}(P \setminus \{\text{event}\})$ is the average correlation between the given recall event's topic vector and
847 the topic vectors from all other recall events not matched to the same video event (for a single
848 participant). We then averaged these distinctiveness scores across all of the events recalled by the
849 given participant and correlated resulting values with hand-annotated and model derived memory
850 performance scores across-subjects, as above.

851 Note: in all instances where we performed statistical tests involving precision or distinctiveness
852 scores, we used Fisher's z-transformation (Fisher, 1925) to stabilize the variance across the distribution
853 of correlation values prior to performing the test. Similarly, when averaging precision or distinctiveness
854 scores, we z-transformed the scores prior to computing the mean, and inverse z-transformed the
855 result.

856 Visualizing the video and recall topic trajectories

857 We used the UMAP algorithm (?) (McInnes et al., 2018) to project the 100-dimensional topic space
858 onto a two-dimensional space for visualization (Figs. 7, 8). To Importantly, to ensure that all of
859 the trajectories were projected onto the *same* lower dimensional space, we computed the low-
860 dimensional embedding on a "stacked" matrix created by vertically concatenating the events-
861 by-topics topic proportions matrices for the video, across-participants average recall and all 17

862 individual participants' recalls. We then divided the rows of the result (a total-number-of-events
863 by two matrix) back into separate matrices for the video topic trajectory, across-participant average
864 recall trajectory and the trajectories for each individual participant's recalls (Fig. 7). This general ap-
865 proach for discovering a shared low-dimensional embedding for a collections of high-dimensional
866 observations follows Heusser et al. (2018b).

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

876 Estimating the consistency of flow through topic space across participants

877 In Figure 7B, we present an analysis aimed at characterizing locations in topic space that dif-
878 ferent participants move through in a consistent way (via their recall topic trajectories). The
879 two-dimensional topic space used in our visualizations (Fig. 7) ~~ranged from -5 to 5 (arbitrary)~~
880 ~~units in the x dimension and from -6.5 to 2 units in the y dimension.~~ We divided this space into a
881 ~~grid of vertices spaced 0.25 units apart~~ comprised a 60 x 60 (arbitrary units) square. We tiled this
882 space with a 50 x 50 grid of evenly spaced vertices, and defined a circular area centered on each
883 vertex whose radius was two times the distance between adjacent vertices (i.e., 2.4 units). For each
884 vertex, we examined the set of line segments formed by connecting each pair successively recalled
885 events, across all participants, that passed ~~within 0.5 units through this circle~~. We computed the
886 distribution of angles formed by those segments and the x-axis, and used a Rayleigh test to deter-
887 mine whether the distribution of angles was reliably "peaked" (i.e., consistent across all transitions
888 that passed through that local portion of topic space). To create Figure 7B we drew an arrow

889 originating from each grid vertex, pointing in the direction of the average angle formed by the line
890 segments that passed within 0.5 units~~its circular radius~~. We set the arrow lengths to be inversely
891 proportional to the p -values of the Rayleigh tests at each vertex. Specifically, for each vertex we
892 converted all of the angles of segments that passed within 0.5~~2.4~~ units to unit vectors, and we set
893 the arrow lengths at each vertex proportional to the length of the (circular) mean vector. We also
894 indicated any significant results ($p < 0.05$, corrected using the Benjamani-Hochberg procedure) by
895 coloring the arrows in blue (darker blue denotes a lower p -value, i.e., a longer mean vector); all
896 tests with $p \geq 0.05$ are displayed in gray and given a lower opacity value.

897 **Searchlight fMRI analyses**

898 In Figure 9, we present two analyses aimed at identifying brain structures~~regions~~ whose responses
899 (as participants viewed the video) exhibited particular temporal correlations~~a particular temporal~~
900 structure. We developed a searchlight analysis whereby~~wherein~~ we constructed a cube~~5 x 5~~
901 x 5 cube of voxels~~(following Chen et al., 2017)~~ centered on each voxel ~~(radius: 5 voxels)~~. For
902 in the brain, and for each of these cubes, ~~we~~ computed the temporal correlation matrix of the
903 voxel responses during video viewing. Specifically, for each of the 1976 volumes collected during
904 video viewing, we correlated the activity patterns in the given cube with the activity patterns (in
905 the same cube) collected during every other timepoint. This yielded a 1976 by 1976 correlation
906 matrix for each cube. Note: participant 5's scan ended 75s early, and in Chen et al., 2017's publicly
907 released dataset, their scan data was padded to match the length of the other participants'. For
908 our searchlight analyses, we removed this padded data (i.e., the last 50 TRs), resulting in a 1925 by
909 1925 correlation matrix for each cube in participant 5's brain.

910 Next, we constructed ~~two sets~~a series of "template" ~~matrices: one reflected the~~" matrices:
911 the first reflecting the timecourse of video's topic trajectory~~and the other reflected~~, and the others
912 reflecting that of each participant's recall topic trajectory. To construct the video template, we
913 computed the correlations between the topic proportions estimated for every pair of TRs (prior to
914 segmenting the trajectory into discrete events; i.e., the correlation matrix shown in Figs. 2B and 9A).
915 We constructed similar temporal correlation matrices for each participant's recall topic trajectory

916 (Figs. 2D, S4). However, to correct for length differences and potential non-linear transformations
917 between viewing time and recall time, we first used dynamic time warping (Berndt and Clifford,
918 1994) to temporally align participants' recall topic trajectories with the video topic trajectory (an-
919 An example correlation matrix before and after warping is shown in Fig. 9B). This yielded a 1976
920 by 1976 correlation matrix for the video template and for each participant's recall template.

921 The temporal structure of the video's content (as described by our model) is captured in the
922 block-diagonal structure of the video's temporal correlation matrix (e.g., Figs. 2B, 9A), with time
923 periods of thematic stability represented as dark blocks of varying sizes. Inspecting the video
924 correlation matrix suggests that the video's semantic content is highly temporally specific (i.e.,
925 the correlations between topic vectors from distant timepoints are almost entirely near-zero).
926 By contrast, the activity patterns of individual (cubes of) voxels can encode relatively limited
927 information on their own, and their activity frequently contributes to multiple separate functions
928 (Freedman et al., 2001; Sigman and Dehaene, 2008; Charron and Koechlin, 2010; Rishel et al., 2013)
929 . By nature, these two attributes give rise to similarities in activity across large timescales that may
930 not necessarily reflect a single task. To enable a more sensitive analysis of brain regions whose shifts
931 in activity patterns mirrored shifts in the semantic content of the video or recalls, we restricted the
932 temporal correlations we considered to timescale of semantic information captured by our model.
933 Specifically, we isolated the upper triangle of the video correlation matrix and created a "proximal
934 correlation mask" that included only diagonals from the upper triangle of the video correlation
935 matrix up to the first that contained no positive correlations. Applying this mask to the full video
936 correlation matrix was analogous to excluding diagonals beyond the corner of the largest diagonal
937 block. In other words, the timescale of temporal correlations we considered corresponded to the
938 longest period of thematic stability in the video, and by extension the longest expected period
939 of thematic stability in participants' recalls and the longest period of stability we might expect
940 to see in voxel activity arising from processing or encoding video content. Figure 9 shows this
941 proximal correlation mask applied to the temporal correlation matrices for the video, an example
942 participant's (warped) recall, and an example cube of voxels from our searchlight analyses.

943 To determine which (cubes of) voxel responses reliably matched the video template, we corre-

944 lated the proximal diagonals from the upper triangle of the voxel correlation matrix for each cube
945 with the upper triangle of the proximal diagonals from video template matrix (Kriegeskorte et al.,
946 2008). This yielded, for each participant, a single correlation value. We computed the average
947 voxelwise map of correlation values. We then performed a one-sample *t*-test on the distribution
948 of (Fisher z-transformed) correlation coefficient correlations at each voxel, across participants. We
949 used This resulted in a value for each voxel (cube), describing how reliably its timecourse mirrored
950 that of the video.

951 We further sought to ensure that our analysis identified regions where the activations' temporal
952 structure specifically reflected that of the video, rather than regions whose activity was simply
953 autocorrelated at a width similar to the video template's diagonal. To achieve this, we used a
954 phase shift-based permutation procedure, wherein we circularly shifted the video's topic trajectory
955 by a permutation-based procedure to assess significance, whereby we re-computed the average
956 correlations for each of 100 "null" videotemplates (constructed by circularly shifting the template
957 by a random number of timepoints), computed the resulting "null" video template, and re-ran
958 the searchlight analysis, in full. (For each permutation of the 100 permutations, the same random
959 shift was used for all participants.) We then z-scored the observed (unshifted) result at each
960 voxel against the distribution of permutation-derived "null" results, and estimated a *p*-value by
961 computing the proportion of shifted correlations that were larger than the observed (unshifted)
962 correlation results that yielded larger values. To create the map in Figure 9A-C, we thresholded out
963 any voxels whose correlation values similarity to the unshifted video's structure fell below the 95th
964 percentile of the permutation-derived null distribution similarity results.

965 We used a similar an analogous procedure to identify which voxels' responses reflected the recall
966 templates. For each participant, we correlated the proximal diagonals from the upper triangle of the
967 correlation matrix for each cube of voxels with their (time-warped the proximal diagonals from the
968 upper triangle of their (time-warped) recall correlation matrix. As in the video template analysis,
969 this yielded a single correlation coefficient for each voxelwise map of correlation coefficients per
970 participant. However, whereas the video analysis compared every participant's responses to
the same template, here the recall templates were unique for each participant. We computed

972 the average As in the analysis described above, we *t*-scored the (Fisher *z*-transformed correlation
973 coefficient across participants) voxelwise correlations, and used the same permutation procedure
974 we developed for the video responses to assess significant correlations ensure specificity to the
975 recall timeseries and assign significance values. To create the map in Figure 9B–we–D we again
976 thresholded out any voxels whose correlation correspondence values fell below the 95th percentile
977 of the permutation-derived null distribution.

978 Neurosynth decoding analyses

979 Neurosynth parses a massive online database of over 14,000 neuroimaging studies and constructs
980 meta-analysis images for over 13,000 psychology- and neuroscience-related terms, based on NIfTI
981 images accompanying studies where those terms appear at a high frequency. Then, given a novel
982 image (tagged with its value type; e.g., *t*-, *F*- or *p*-statistics), Neurosynth returns a list of terms whose
983 meta-analysis images are most similar to this new data. Our permutation procedure yielded, for
984 each of the two searchlight analyses, a voxelwise map of significance (*p*-statistic) values. These
985 maps describe the extent to which each voxel specifically reflected the temporal structure of the video
986 or individuals' recalls (i.e., for each voxel, the proportion of phase-shifted topic vector correlation
987 matrices less similar to the voxel activity correlation matrix than the unshifted video's correlation
988 matrix). We input the two statistical maps described above to Neurosynth to create a list of the 10
989 most representative terms for each map.

990 **References**

- 991 Atkinson, R. C. and Shiffrin, R. M. (1968). Human memory: A proposed system and its control
992 processes. In Spence, K. W. and Spence, J. T., editors, *The psychology of learning and motivation*,
993 volume 2, pages 89–105. Academic Press, New York.
- 994 Baldassano, C., Chen, J., Zadbood, A., Pillow, J. W., Hasson, U., and Norman, K. A. (2017).

- 995 Discovering event structure in continuous narrative perception and memory. *Neuron*, 95(3):709–
996 721.
- 997 Baldassano, C., Hasson, U., and Norman, K. A. (2018). Representation of real-world event schemas
998 during narrative perception. *Journal of Neuroscience*, 38(45):9689–9699.
- 999 Berndt, D. J. and Clifford, J. (1994). Using dynamic time warping to find patterns in time series. In
1000 *KDD workshop*, volume 10, pages 359–370.
- 1001 Blei, D. M. and Lafferty, J. D. (2006). Dynamic topic models. In *Proceedings of the 23rd International*
1002 *Conference on Machine Learning*, ICML '06, pages 113–120, New York, NY, US. ACM.
- 1003 Blei, D. M., Ng, A. Y., and Jordan, M. I. (2003). Latent dirichlet allocation. *Journal of Machine*
1004 *Learning Research*, 3:993 – 1022.
- 1005 Brunec, I. K., Moscovitch, M. M., and Barense, M. D. (2018). Boundaries shape cognitive represen-
1006 tations of spaces and events. *Trends in Cognitive Sciences*, 22(7):637–650.
- 1007 Bylinskii, Z., Isola, P., Bainbridge, C., Torralba, A., and Oliva, A. (2015). Intrinsic and extrinsic
1008 effects on image memorability. *Vision Research*, 116:165–178.
- 1009 Capota, M., Turek, J., Chen, P.-H., Zhu, X., Manning, J. R., Sundaram, N., Keller, B., Wang, Y., and
1010 Shin, Y. S. (2017). Brain imaging analysis kit.
- 1011 Cer, D., Yang, Y., Kong, S. Y., Hua, N., Limtiaco, N., John, R. S., Constant, N., Guajardo-Cespedes,
1012 M., Yuan, S., Tar, C., Sung, Y.-H., Strope, B., and Kurzweil, R. (2018). Universal sentence encoder.
1013 *arXiv*, 1803.11175.
- 1014 Charron, S. and Koechlin, E. (2010). Divided representations of current goals in the human frontal
1015 lobes. *Science*, 328(5976):360–363.
- 1016 Chen, J., Leong, Y. C., Honey, C. J., Yong, C. H., Norman, K. A., and Hasson, U. (2017). Shared
1017 memories reveal shared structure in neural activity across individuals. *Nature Neuroscience*,
1018 20(1):115.

- 1019 Chun, M. and Turk-Browne, N. (2007). Interactions between attention and memory. *Current opinion*
1020 *in neurobiology*, 17(2):177–184.
- 1021 Clewett, D. and Davachi, L. (2017). The ebb and flow of experience determines the temporal
1022 structure of memory. *Curr Opin Behav Sci*, 17:186–193.
- 1023 Davachi, L. (2006). Item, context and relational episodic encoding in humans. *Current Opinion in*
1024 *Neurobiology*, 16(6):693—700.
- 1025 Davachi, L., Mitchell, J. P., and Wagner, A. D. (2003). Multiple routes to memory: distinct medial
1026 temporal lobe processes build item and source memories. *Proceedings of the National Academy of*
1027 *Sciences, USA*, 100(4):2157 – 2162.
- 1028 Dobrushin, R. L. (1970). Prescribing a system of random variables by conditional distributions.
1029 *Theory of Probability & Its Applications*, 15(3):458–486.
- 1030 DuBrow, S. and Davachi, L. (2013). The influence of contextual boundaries on memory for the
1031 sequential order of events. *Journal of Experimental Psychology: General*, 142(4):1277–1286.
- 1032 Ezzyat, Y. and Davachi, L. (2011). What constitutes an episode in episodic memory? *Psychological*
1033 *Science*, 22(2):243–252.
- 1034 Fisher, R. A. (1925). *Statistical Methods for Research Workers*. Oliver and Boyd.
- 1035 Freedman, D., Riesenhuber, M., Poggio, T., and Miller, E. (2001). Categorical representation of
1036 visual stimuli in the primate prefrontal cortex. *Science*, 291(5502):312–316.
- 1037 Friendly, M., Franklin, P. E., Hoffman, D., and Rubin, D. C. (1982). The Toronto Word Pool:
1038 Norms for imagery, concreteness, orthographic variables, and grammatical usage for 1,080
1039 words. *Behavior Research Methods and Instrumentation*, 14:375–399.
- 1040 Gilboa, A. and Marlatte, H. (2017). Neurobiology of schemas and schema-mediated memory.
1041 *Trends Cogn Sci*, 21(8):618–631.

- 1042 Hasson, U., Chen, J., and Honey, C. J. (2015). Hierarchical process memory: memory as an integral
1043 component of information processing. *Trends in Cognitive Science*, 19(6):304–315.
- 1044 Hasson, U., Yang, E., Vallines, I., Heeger, D. J., and Rubin, N. (2008). A hierarchy of temporal
1045 receptive windows in human cortex. *Journal of Neuroscience*, 28(10):2539–2550.
- 1046 Heusser, A. C., Ezzyat, Y., Shiff, I., and Davachi, L. (2018a). Perceptual boundaries cause mnemonic
1047 trade-offs between local boundary processing and across-trial associative binding. *Journal of
1048 Experimental Psychology Learning, Memory, and Cognition*, 44(7):1075–1090.
- 1049 Heusser, A. C., Fitzpatrick, P. C., Field, C. E., Ziman, K., and Manning, J. R. (2017). Quail: a
1050 Python toolbox for analyzing and plotting free recall data. *The Journal of Open Source Software*,
1051 10.21105/joss.00424.
- 1052 Heusser, A. C., Ziman, K., Owen, L. L. W., and Manning, J. R. (2018b). HyperTools: a Python
1053 toolbox for gaining geometric insights into high-dimensional data. *Journal of Machine Learning
1054 Research*, 18(152):1–6.
- 1055 Howard, M. W. and Kahana, M. J. (2002). A distributed representation of temporal context. *Journal
1056 of Mathematical Psychology*, 46:269–299.
- 1057 Howard, M. W., MacDonald, C. J., Tiganj, Z., Shankar, K. H., Du, Q., Hasselmo, M. E., and H., E.
1058 (2014). A unified mathematical framework for coding time, space, and sequences in the medial
1059 temporal lobe. *Journal of Neuroscience*, 34(13):4692–4707.
- 1060 Howard, M. W., Viskontas, I. V., Shankar, K. H., and Fried, I. (2012). Ensembles of human MTL
1061 neurons “jump back in time” in response to a repeated stimulus. *Hippocampus*, 22:1833–1847.
- 1062 Huk, A., Bonnen, K., and He, B. J. (2018). Beyond trial-based paradigms: continuous behavior, on-
1063 going neural activity, and naturalistic stimuli. *Journal of Neuroscience*, 10.1523/JNEUROSCI.1920-
1064 17.2018.
- 1065 Kahana, M. J. (1996). Associative retrieval processes in free recall. *Memory & Cognition*, 24:103–109.

- 1066 Kahana, M. J. (2012). *Foundations of Human Memory*. Oxford University Press, New York, NY.
- 1067 Koriat, A. and Goldsmith, M. (1994). Memory in naturalistic and laboratory contexts: distin-
1068 guishing accuracy-oriented and quantity-oriented approaches to memory assessment. *Journal of*
1069 *Experimental Psychology: General*, 123(3):297–315.
- 1070 Kriegeskorte, N., Mur, M., and Bandettini, P. (2008). Representational similarity analysis – con-
1071 nnecting the branches of systems neuroscience. *Frontiers in Systems Neuroscience*, 2:1 – 28.
- 1072 Landauer, T. K., Foltz, P. W., and Laham, D. (1998). Introduction to latent semantic analysis.
1073 *Discourse Processes*, 25:259–284.
- 1074 Lerner, Y., Honey, C. J., Silbert, L. J., and Hasson, U. (2011). Topographic mapping of a hierarchy
1075 of temporal receptive windows using a narrated story. *Journal of Neuroscience*, 31(8):2906–2915.
- 1076 Manning, J. R. (2019). Episodic memory: mental time travel or a quantum ‘memory wave’ function?
1077 *PsyArXiv*, doi:10.31234/osf.io/6zjwb.
- 1078 Manning, J. R., Norman, K. A., and Kahana, M. J. (2015). The role of context in episodic memory.
1079 In Gazzaniga, M., editor, *The Cognitive Neurosciences, Fifth edition*, pages 557–566. MIT Press.
- 1080 Manning, J. R., Polyn, S. M., Baltuch, G., Litt, B., and Kahana, M. J. (2011). Oscillatory patterns
1081 in temporal lobe reveal context reinstatement during memory search. *Proceedings of the National*
1082 *Academy of Sciences, USA*, 108(31):12893–12897.
- 1083 McInnes, L., Healy, J., and Melville, J. (2018). UMAP: Uniform manifold approximation and
1084 projection for dimension reduction. *arXiv*, 1802(03426).
- 1085 Menon, V. and Uddin, L. Q. (2010). Saliency, switching, attention and control: a network model of
1086 insula function. *Brain Structure and Function*, 214(5-6):655–667.
- 1087 Mikolov, T., Chen, K., Corrado, G., and Dean, J. (2013). Efficient estimation of word representations
1088 in vector space. *arXiv*, 1301.3781.

- 1089 Mueller, A., Fillion-Robin, J.-C., Boidol, R., Tian, F., Nechifor, P., yoonsubKim, Peter, Rampin, R.,
1090 Corvellec, M., Medina, J., Dai, Y., Petrushev, B., Langner, K. M., Hong, Alessio, Ozsvald, I.,
1091 vkolmakov, Jones, T., Bailey, E., Rho, V., IgorAPM, Roy, D., May, C., foobuzz, Piyush, Seong,
1092 L. K., Goey, J. V., Smith, J. S., Gus, and Mai, F. (2018). WordCloud 1.5.0: a little word cloud
1093 generator in Python. *Zenodo*, <https://zenodo.org/record/1322068#.W4tPKZNKh24>.
- 1094 Murdock, B. B. (1962). The serial position effect of free recall. *Journal of Experimental Psychology*,
1095 64:482–488.
- 1096 Paller, K. A. and Wagner, A. D. (2002). Observing the transformation of experience into memory.
1097 *Trends in Cognitive Sciences*, 6(2):93–102.
- 1098 Pedregosa, F., Varoquaux, G., Gramfort, A., Michel, V., Thirion, B., Grisel, O., Blondel, M., Pretten-
1099 hofer, P., Weiss, R., Dubourg, V., Vanderplas, J., Passos, A., Cournapeau, D., Brucher, M., Perrot,
1100 M., and Duchesnay, E. (2011). Scikit-learn: Machine learning in Python. *Journal of Machine
1101 Learning Research*, 12:2825–2830.
- 1102 Polyn, S. M., Norman, K. A., and Kahana, M. J. (2009). A context maintenance and retrieval model
1103 of organizational processes in free recall. *Psychological Review*, 116(1):129–156.
- 1104 Postman, L. and Phillips, L. W. (1965). Short-term temporal changes in free recall. *Quarterly Journal
1105 of Experimental Psychology*, 17:132–138.
- 1106 Rabiner, L. (1989). A tutorial on Hidden Markov Models and selected applications in speech
1107 recognition. *Proceedings of the IEEE*, 77(2):257–286.
- 1108 Radvansky, G. A. and Zacks, J. M. (2017). Event boundaries in memory and cognition. *Curr Opin
1109 Behav Sci*, 17:133–140.
- 1110 Ramdas, A., Trillos, N., and Cuturi, M. (2017). On wasserstein two-sample testing and related
1111 families of nonparametric tests. *Entropy*, 19(2):47.
- 1112 Ranganath, C. and Ritchey, M. (2012). Two cortical systems for memory-guided behavior. *Nature
1113 Reviews Neuroscience*, 13:713 – 726.

- 1114 Rishel, C. A., Huang, G., and Freedman, D. J. (2013). Independent category and spatial encoding
1115 in parietal cortex. *Neuron*, 77(5):969–979.
- 1116 Schlichting, M. L. and Preston, A. R. (2015). Memory integration: neural mechanisms and impli-
1117 cations for behavior. *Current Opinion in Behavioral Sciences*, 1:1–8.
- 1118 Sigman, M. and Dehaene, S. (2008). Brain mechanisms of serial and parallel processing during
1119 dual-task performance. *Journal of Neuroscience*, 28(30):7585–7589.
- 1120 Simony, E., Honey, C. J., Chen, J., and Hasson, U. (2016). Uncovering stimulus-locked network
1121 dynamics during narrative comprehension. *Nature Communications*, 7(12141):1–13.
- 1122 Spalding, K. N., Schlichting, M. L., Zeithamova, D., Preston, A. R., Tranel, D., Duff, M. C., and
1123 Warren, D. E. (2018). Ventromedial prefrontal cortex is necessary for normal associative inference
1124 and memory integration. *The Journal of Neuroscience*, 38(15):3767–3775.
- 1125 Stark, S. M., Yassa, M. A., and Stark, C. E. L. (2010). Individual differences in spatial pattern
1126 separation performance associated with healthy aging in humans. *Learning & Memory*, 17(6):284–
1127 288.
- 1128 Steyvers, M., Shiffrin, R. M., and Nelson, D. L. (2004). Word association spaces for predicting
1129 semantic similarity effects in episodic memory. In Healy, A. F., editor, *Cognitive Psychology and*
1130 *its Applications: Festschrift in Honor of Lyle Bourne, Walter Kintsch, and Thomas Landauer*. American
1131 Psychological Association, Washington, DC.
- 1132 Tompry, A. and Davachi, L. (2017). Consolidation promotes the emergence of representational
1133 overlap in the hippocampus and medial prefrontal cortex. *Neuron*, 96(1):228–241.
- 1134 Tyng, C. M., Amin, H. U., Saad, M. N. M., and S, M. A. (2017). The influences of emotion on
1135 learning and memory. *Frontiers in psychology*, 8:1454.
- 1136 van Kesteren, M. T. R., Ruiter, D. J., Fernández, G., and Henson, R. N. (2012). How schema and
1137 novelty augment memory formation. *Trends Neurosci*, 35(4):211–9.

- 1138 Waskom, M., Botvinnik, O., Okane, D., Hobson, P., David, Halchenko, Y., Lukauskas, S., Cole, J. B.,
1139 Warmenhoven, J., de Ruiter, J., Hoyer, S., Vanderplas, J., Villalba, S., Kunter, G., Quintero, E.,
1140 Martin, M., Miles, A., Meyer, K., Augspurger, T., Yarkoni, T., Bachant, P., Williams, M., Evans,
1141 C., Fitzgerald, C., Brian, Wehner, D., Hitz, G., Ziegler, E., Qalieh, A., and Lee, A. (2016). Seaborn:
1142 v0.7.1.
- 1143 Welch, G. B. and Burnett, C. T. (1924). Is primacy a factor in association-formation. *American Journal
1144 of Psychology*, 35:396–401.
- 1145 Yarkoni, T., Poldrack, R. A., Nichols, T. E., Van Essen, D. C., and Wager, T. D. (2011). Large-scale
1146 automated synthesis of human functional neuroimaging data. *Nature Methods*, 8(8):665.
- 1147 Yassa, M. A., Lacy, J. W., Stark, S. M., Albert, M. S., Gallagher, M., and Stark, C. E. L. (2011). Pattern
1148 separation deficits associated with increased hippocampal ca3 and dentate gyrus activity in
1149 nondemented older adults. *Hippocampus*, 21(9):968–979.
- 1150 Yassa, M. A. and Stark, C. E. L. (2011). Pattern separation in the hippocampus. *Trends In Neuro-
1151 sciences*, 34(10):515–525.
- 1152 Yonelinas, A. P. (2002). The nature of recollection and familiarity: A review of 30 years of research.
1153 *Journal of Memory and Language*, 46:441–517.
- 1154 Yonelinas, A. P., Kroll, N. E., Quamme, J. R., Lazzara, M. M., Sauvé, M. J., Widaman, K. F., and
1155 Knight, R. T. (2002). Effects of extensive temporal lobe damage or mild hypoxia on recollection
1156 and familiarity. *Nature Neuroscience*, 5(11):1236–41.
- 1157 Zacks, J. M., Speer, N. K., Swallow, K. M., Braver, T. S., and Reynolds, J. R. (2007). Event perception:
1158 a mind-brain perspective. *Psychological Bulletin*, 133:273–293.
- 1159 Zadbood, A., Chen, J., Leong, Y. C., Norman, K. A., and Hasson, U. (2017). How we transmit
1160 memories to other brains: Constructing shared neural representations via communication. *Cereb
1161 Cortex*, 27(10):4988–5000.

¹¹⁶² Zwaan, R. A. and Radvansky, G. A. (1998). Situation models in language comprehension and
¹¹⁶³ memory. *Psychological Bulletin*, 123(2):162 – 185.

¹¹⁶⁴ Supporting information

¹¹⁶⁵ Supporting information is available in the online version of the paper.

¹¹⁶⁶ Acknowledgements

¹¹⁶⁷ We thank Luke Chang, Janice Chen, Chris Honey, Lucy Owen, Emily Whitaker, and Kirsten Ziman
¹¹⁶⁸ for feedback and scientific discussions. We also thank Janice Chen, Yuan Chang Leong, Kenneth
¹¹⁶⁹ Norman, and Uri Hasson for sharing the data used in our study. Our work was supported in part
¹¹⁷⁰ by NSF EPSCoR Award Number 1632738. The content is solely the responsibility of the authors
¹¹⁷¹ and does not necessarily represent the official views of our supporting organizations.

¹¹⁷² Author contributions

¹¹⁷³ Conceptualization: A.C.H. and J.R.M.; Methodology: A.C.H., P.C.F. and J.R.M.; Software: A.C.H.,
¹¹⁷⁴ 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.,
¹¹⁷⁵ P.C.F. and J.R.M.; Supervision: J.R.M.

¹¹⁷⁶ Author information

¹¹⁷⁷ The authors declare no competing financial interests. Correspondence and requests for materials
¹¹⁷⁸ should be addressed to J.R.M. (jeremy.r.manning@dartmouth.edu).