

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

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

Department of Psychological and Brain Sciences

Dartmouth College, Hanover, NH 03755, USA

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

⁵ July 23, 2019

⁶ **Abstract**

⁷ Our life experiences unfold over time in highly complex manner, with the evolving presence
⁸ and absence of numerous intricate features describing our journey between each circumstance
⁹ or event we encounter. Here, we propose a framework for mapping dynamic naturalistic ex-
¹⁰ periences onto geometric spaces as *trajectories* that capture the temporal dynamics of real-world
¹¹ content. Within this geometric framework, one may compare the shape of the trajectory formed
¹² by an experience to that defined by one's later recollection to characterize our memories' re-
¹³ covery and distortion of the external world. Here, we apply this approach to a naturalistic
¹⁴ memory experiment in which participants viewed and verbally recounted a video, and find
¹⁵ that the video and subsequent recalls share both an experience-specific shape and a discernible
¹⁶ event-like structure. However, the level of *precision* with which individuals recounted various
¹⁷ events and the *distinctiveness* of recall for those events were varied and predictive of overall
¹⁸ memory performance. Finally, we identify a network of brain structures that is sensitive to the
¹⁹ "shapes" of our ongoing experiences, and an overlapping network sensitive to how we will later

remember them. These results highlight the rich event-like structure of the external world and our memories, and offer novel, content-sensitive alternatives to classical “proportion recalled” measures for assessing episodic memory

Introduction

What does it mean to *remember* something? In traditional episodic memory experiments (e.g., list-learning or trial-based experiments; Murdock, 1962; Kahana, 1996), remembering is often cast as a discrete and binary operation: each studied item may be separated from the rest of one’s experiences, and that item may be labeled as having been recalled versus forgotten. More nuanced studies might incorporate self-reported confidence measures as a proxy for memory strength, or ask participants to discriminate between “recollecting” the (contextual) details of an experience or having a general feeling of “familiarity” (Yonelinas, 2002). Using well-controlled, trial-based experimental designs, the field has amassed a wealth of valuable information regarding human episodic memory. However, there are fundamental properties of the external world and our memories that trial-based experiments are not well suited to capture (for review also see Kriat and Goldsmith, 1994; Huk et al., 2018). First, our experiences and memories are continuous, rather than discrete—removing a (naturalistic) event from the context in which it occurs can substantially change its meaning. Second, the specific language used to describe an experience has little bearing on whether the experience 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, the number or proportion of precise recalls is often a primary metric for assessing the quality of participants’ memories. Third, one might remember the *essence* (or a general summary) of an experience but forget (or neglect to recount) particular details. Capturing the essence of what happened is typically the main “point” of recounting a memory to a listener, while the addition of highly specific details may add comparatively little to successful conveyance of an experience.

46 How might one go about formally characterizing the essence of an experience, or whether it
47 has been recovered by the rememberer? Any given moment of an experience derives meaning
48 from surrounding moments, as well as from longer-range temporal associations (e.g., Lerner et al.,
49 2011; Manning, 2019). Therefore, the timecourse describing how an event unfolds is fundamental
50 to its overall meaning. Further, this hierarchy formed by our subjective experiences at different
51 timescales defines a *context* for each new moment (e.g., Howard and Kahana, 2002; Howard et al.,
52 2014), and plays an important role in how we interpret that moment and remember it later (for
53 review see Manning et al., 2015). Our memory systems can leverage these associations to form
54 predictions that help guide our behaviors (Ranganath and Ritchey, 2012). For example, as we
55 navigate the world, the features of our subjective experiences tend to change gradually (e.g., the
56 room or situation we are in at any given moment is strongly temporally autocorrelated), allowing
57 us to form stable estimates of our current situation and behave accordingly (Zacks et al., 2007;
58 Zwaan and Radvansky, 1998).

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

72 Here we sought to examine how the temporal dynamics of a “naturalistic” experience were
73 later reflected in participants’ memories. We analyzed an open dataset that comprised behavioral

74 and functional Magnetic Resonance Imaging (fMRI) data collected as participants viewed and then
75 verbally recounted an episode of the BBC television series *Sherlock* (Chen et al., 2017). We developed
76 a computational framework for characterizing the temporal dynamics of the moment-by-moment
77 content of the episode and of participants' verbal recalls. Specifically, we use topic modeling (Blei
78 et al., 2003) to characterize the thematic conceptual (semantic) content present in each moment of
79 the episode and recalls, and Hidden Markov Models (Rabiner, 1989; Baldassano et al., 2017) to
80 discretize the evolving semantic content into events. In this way, we cast naturalistic experiences
81 (and recalls of those experiences) as *trajectories* that describe how the experiences evolve over
82 time. Under this framework, successful remembering entails verbally "traversing" the content
83 trajectory of the episode, thereby reproducing the shape (or essence) of the original experience.
84 Comparing the shapes of the topic trajectories of the episode and of participants' retellings of the
85 episode then reveals which aspects of the episode were preserved (or lost) in the translation into
86 memory. We further examine whether 1) the *precision* with which a participant recounts each event
87 and 2) the *distinctiveness* each recall event is (relative to the other recalled events) relates to their
88 overall memory performance. Last, we identify networks of brain structures whose responses (as
89 participants watched the episode) reflected the shape of the episode, and how participants would
90 later recount the episode.

91 **Results**

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

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

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

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



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

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

145 Given that the time-varying content of the video could be segmented cleanly into discrete
146 events, we wondered whether participants’ recalls of the video also displayed a similar structure.
147 We applied the same topic model (already trained on the video annotations) to each participant’s
148 recalls. Analogous to how we analyzed the time-varying content of the video, to obtain similar
149 estimates for participants’ recalls, we treated each (overlapping) 10 sentence “window” of their
150 transcript as a “document” and then computed the most probable mix of topics reflected in each
151 timepoint’s sentences. This yielded, for each participant, a number-of-sentences by number-of-
152 topics topic proportions matrix that characterized how the topics identified in the original video
153 were reflected in the participant’s recalls. Note that an important feature of our approach is
154 that it allows us to compare participant’s recalls to events from the original video, despite that
155 different participants may have used different language to describe the same event, and that those

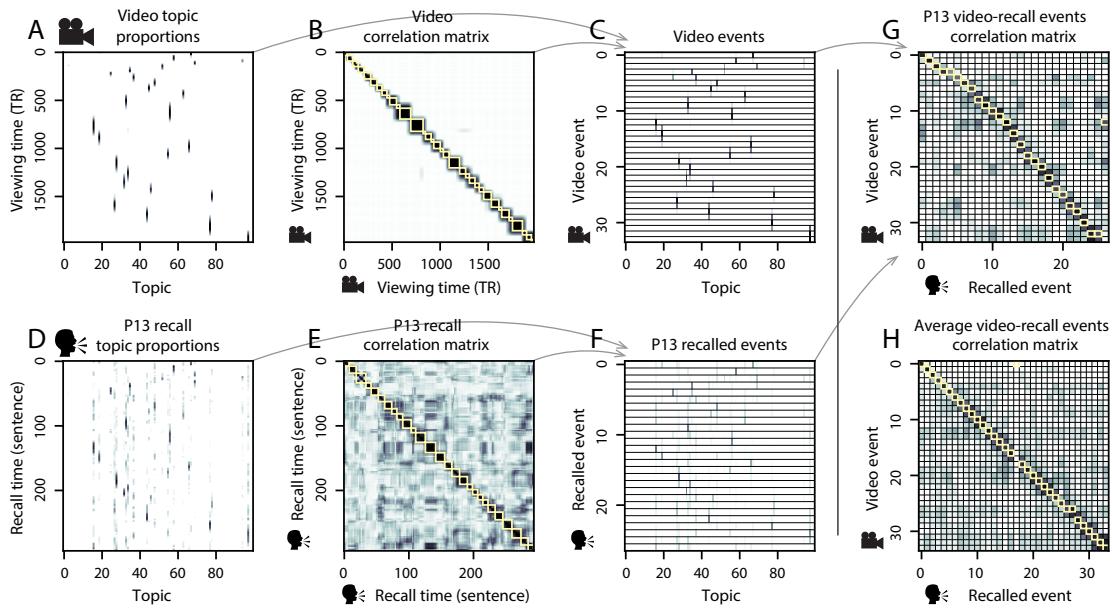


Figure 2: Modelling naturalistic stimuli and recalls. All panels: darker colors indicate greater values; range: [0, 1]. **A.** Topic vectors ($K = 100$) for each of the 1976 video timepoints. **B.** Timepoint-by-timepoint correlation matrix of the topic vectors displayed in Panel A. Event boundaries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

descriptions may not match the original annotations. This is a substantial benefit of projecting the video and recalls into a shared “topic” space. An example topic proportions matrix from one participant’s recalls is shown in Figure 2D.

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

Two clear patterns emerged from this set of analyses. First, although every individual participant’s recalls could be segmented into discrete events (i.e., every individual participant’s recall correlation matrix exhibited clear block diagonal structure; Fig. S4), each participant appeared to have a unique *recall resolution*, reflected in the sizes of those blocks. For example, some participants’ recall topic proportions segmented into just a few events (e.g., Participants P1, P4, and P15), while others’ recalls segmented into many shorter duration events (e.g., Participants P12, P13, and P17). This suggests that different participants may be recalling the video with different levels of detail—e.g., some might touch on just the major plot points, whereas others might attempt to recall every minor scene or action. The second clear pattern present in every individual participant’s recall correlation matrix is that, unlike in the video correlation matrix, there are substantial off-diagonal correlations in participant’s recalls. Whereas each event in the original video was (largely) separable from the others (Fig. 2B), in transforming those separable events into memory participants

¹⁸⁴ appear to be integrating *across* different events, blending elements of previously recalled and not-
¹⁸⁵ yet-recalled events into each newly recalled event (Figs. 2D, S4; also see Manning et al., 2011;
¹⁸⁶ Howard et al., 2012).

¹⁸⁷ The above results indicate that both the structure of the original video and participants' recalls
¹⁸⁸ of the video exhibit event boundaries that can be identified automatically by characterizing the
¹⁸⁹ dynamic content using a shared topic model and segmenting the content into events using HMMs.
¹⁹⁰ Next, we asked whether some correspondence might be made between the specific content of the
¹⁹¹ events the participants experienced in the video, and the events they later recalled. One approach
¹⁹² to linking the experienced (video) and recalled events is to label each recalled event as matching the
¹⁹³ video event with the most similar (i.e., most highly correlated) topic vector (Figs. 2G, S5). This yields
¹⁹⁴ a sequence of "presented" events from the original video, and a sequence of (potentially differently
¹⁹⁵ ordered) "recalled" events for each participant. Analogous to classic list-learning studies, we
¹⁹⁶ can then examine participants' recall sequences by asking which events they tended to recall
¹⁹⁷ first (probability of first recall; Fig. 3A; Welch and Burnett, 1924; Postman and Phillips, 1965;
¹⁹⁸ Atkinson and Shiffrin, 1968); how participants most often transition between recalls of the events
¹⁹⁹ as a function of the temporal distance between them (lag-conditional response probability; Fig. 3B;
²⁰⁰ Kahana, 1996); and which events they were likely to remember overall (serial position recall
²⁰¹ analyses; Fig. 3C; Murdock, 1962). Interestingly, for two of these analyses (probability of first recall
²⁰² and lag-conditional response probability curves) we observe patterns comparable to classic effects
²⁰³ from the list-learning literature: namely, a higher probability of initiating recall with the first event
²⁰⁴ in the sequence (Fig. 3A) and a higher probability of transitioning to neighboring events with an
²⁰⁵ asymmetric forward bias (Fig. 3C). In contrast, we do not observe a pattern comparable to the serial
²⁰⁶ position effect (Fig. 3C), but rather we see higher memory for specific events distributed somewhat
²⁰⁷ evenly throughout the video.

²⁰⁸ Statistical models of memory studies often treat memory recalls as binary (e.g. the item was re-
²⁰⁹ called or not) and independent events. However, our framework produces a content-based model
²¹⁰ of individual stimulus and recall events, allowing for direct quantitative comparison between all
²¹¹ stimulus and recall events, as well as between the recall events themselves. Leveraging these

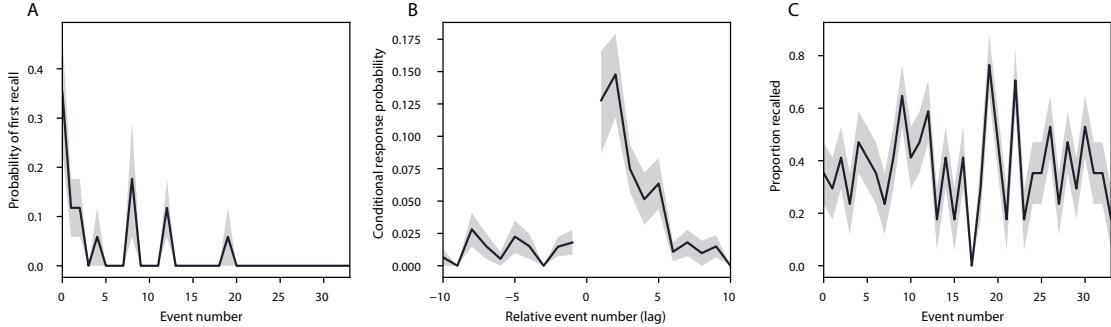


Figure 3: Naturalistic extensions of classic list-learning memory analyses. A. The probability of first recall as a function of the serial position of the event in the video. B. The probability of recalling each event, conditioned on having most recently recalled the event *lag* events away in the video.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

content-based models of the stimulus/recall events, we developed two novel metrics for quantifying naturalistic memory representations: *precision* and *distinctiveness*. We define precision as the average correlation between the topic proportions of each recall event and the maximally correlated video event (Fig. 4). Participants whose recall events are more veridical descriptions of what happened in the video event will presumably have higher precision scores. We find that across participants, a higher precision score is correlated to both hand annotated memory performance (Pearson's $r(15) = 0.6, p = 0.011$) and the number of recall events estimated by our model (Pearson's $r(15) = 0.64, p = 0.005$). A second novel metric we introduce here is distinctiveness, or how unique the recall description was to each video event. We define distinctiveness as 1 minus the average of all non-matching recall events from the video-recall correlation matrix. We hypothesized that participants who recounted events in a more distinctive way would display better overall memory. Similarly to precision, we find that the more distinct participants recalls are (on average), the more they remembered (hand-annotated memory: Pearson's $r(15) = 0.83, p < 0.001$; model-derived memory: Pearson's $r(15) = 0.71, p = 0.001$). In summary, using two novel metrics afforded by our approach, we find that participants whose recalls are both more precise and distinct remember more content.

The prior analyses leverage the correspondence between the 100-dimensional topic proportion

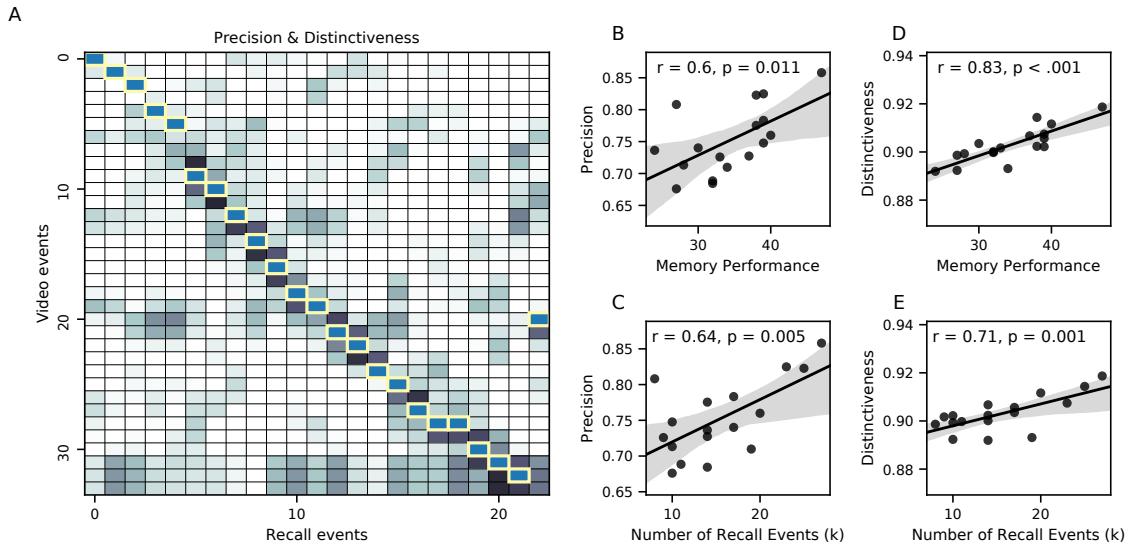


Figure 4: Novel content-based metrics of naturalistic memory: precision and distinctiveness. **A.** A video-recall correlation matrix for a representative participant (17). The yellow boxes highlight the maximum correlation in each column. Precision was computed as the average of the maximum correlation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

229 matrices for the video and participants' recalls to characterize recall. However, it is difficult
230 to gain deep insights into that content solely by examining the topic proportion matrices (e.g.,
231 Figs. 2A, D) or the corresponding correlation matrices (Figs. 2B, E, S4). To visualize the time-
232 varying high-dimensional content in a more intuitive way (Heusser et al., 2018b) we projected the
233 topic proportions matrices onto a two-dimensional space using Uniform Manifold Approximation
234 and Projection (UMAP; McInnes and Healy, 2018). In this lower-dimensional space, each point
235 represents a single video or recall event, and the distances between the points reflect the distances
236 between the events' associated topic vectors (Fig. 5). In other words, events that are near to each
237 other in this space are more semantically similar.

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

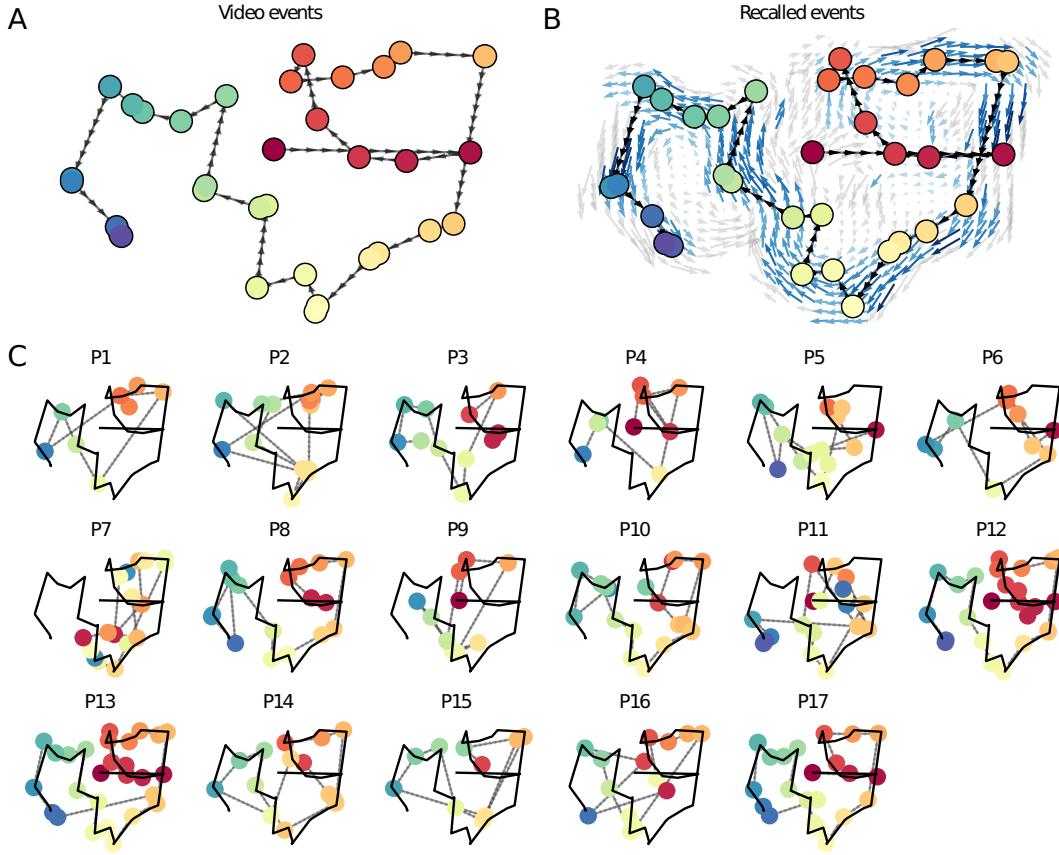


Figure 5: Trajectories through topic space capture the dynamic content of the video and recalls. All panels: the topic proportion matrices have been projected onto a shared two-dimensional space using UMAP. **A.** The two-dimensional topic trajectory taken by the episode of *Sherlock*. Each dot indicates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

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

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

284 In addition to constructing overall summaries, assessing the video and recall topic vectors

from individual events can provide further insights. Specifically, for any given event we can construct two wordles: one from the original video event's topic vector, and a second from the average topic vectors produced by all participants' recalls of that event. We can then examine those wordles visually to gain an intuition for which aspects of the video event were recapitulated in participants' recalls. Several example wordles are displayed in Figure 6C (wordles from the three best-remembered events are circled in green; wordles from the three worst-remembered events are circled in red). Using wordles to visually compare the topical content of each video event and the (average) corresponding recall event reveals the specific content from the specific events that is reliably retained in the transformation into memory (green events) or not (red events).

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

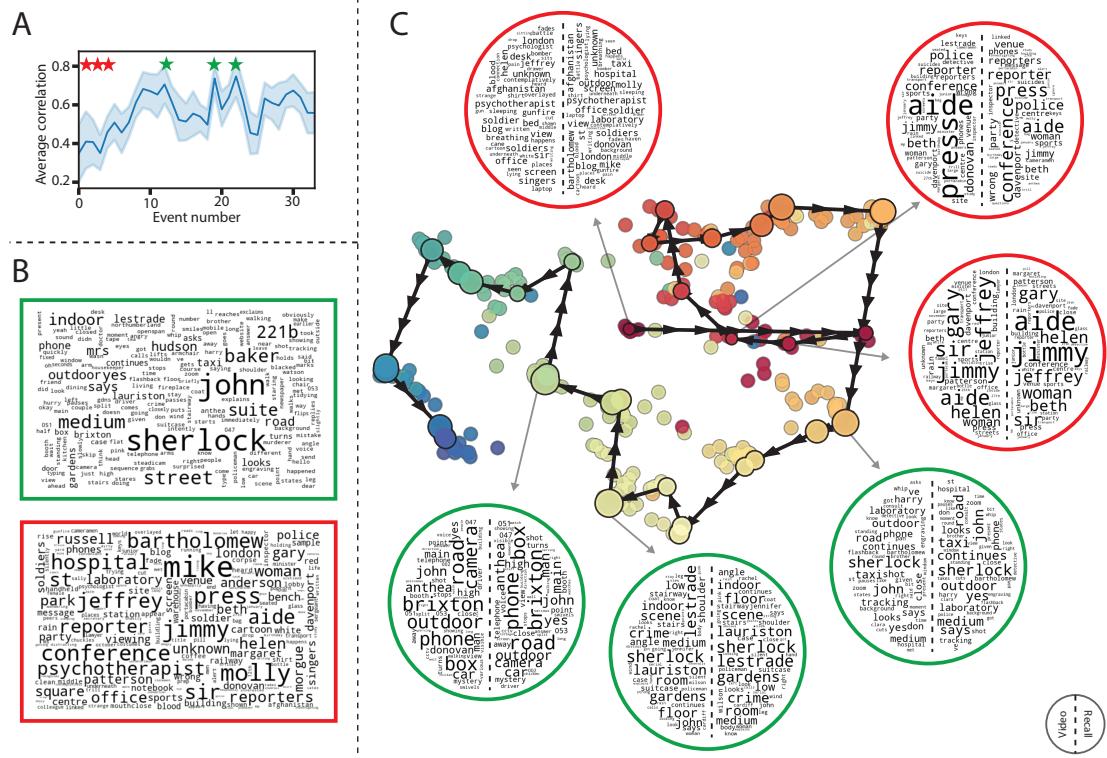


Figure 6: Transforming experience into memory. **A.** Average correlations (across participants) between the topic vectors from each video event and the closest-matching recall events. Error bars denote bootstrap-derived across-participant 95% confidence intervals. The stars denote the three best-remembered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

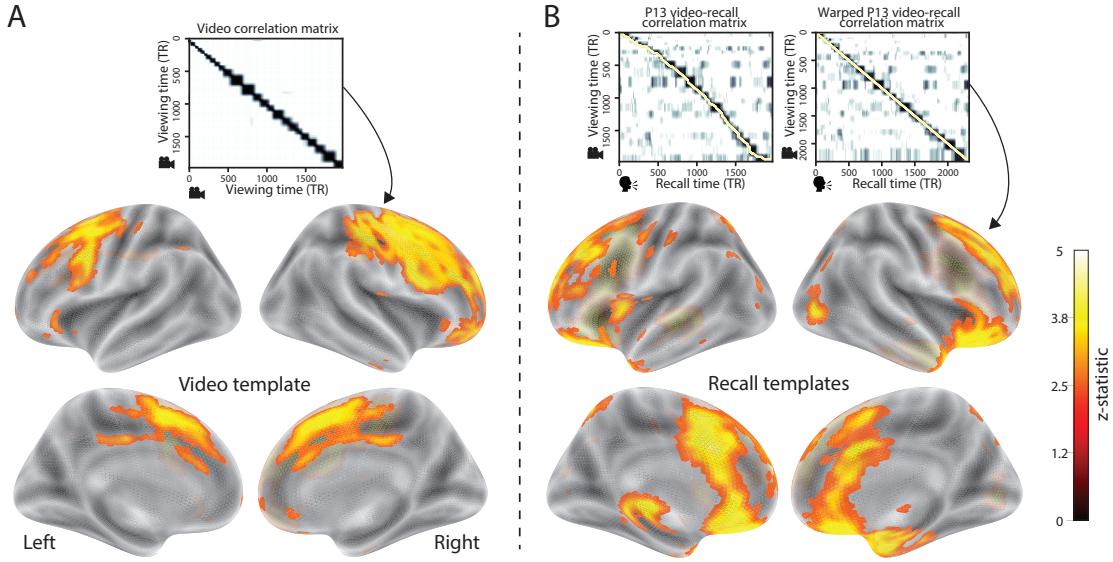


Figure 7: Brain structures that underlie the transformation of experience into memory. **A.** We searched for regions whose responses (as participants watched the video) matched the temporal correlation matrix of the video topic proportions. These regions are sensitive to the narrative structure of the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

classic subsequent memory analyses (e.g., Paller and Wagner, 2002) to domain of naturalistic stimuli.

Discussion

Our work casts remembering as reproducing (behaviorally and neurally) the topic trajectory, or shape, of an experience. This view draws inspiration from prior work aimed at elucidating the neural and behavioral underpinnings of how we process dynamic naturalistic experiences and remember them later. One approach to identifying neural responses to naturalistic stimuli (including experiences) entails building a model of the stimulus and searching for brain regions whose responses are consistent with the model. In prior work, a series of studies from Uri Hasson's group (Lerner et al., 2011; Simony et al., 2016; Chen et al., 2017; Baldassano et al., 2017; Zadbood et al., 2017) have extended this approach with a clever twist: rather than building an

323 explicit stimulus model, these studies instead search for brain responses (while experiencing the
324 stimulus) that are reliably similar across individuals. So called *inter-subject correlation* (ISC) and
325 *inter-subject functional connectivity* (ISFC) analyses effectively treat other people's brain responses
326 to the stimulus as a "model" of how its features change over time. By contrast, in our present
327 work we used topic models and HMMs to construct an explicit stimulus model (i.e., the topic
328 trajectory of the video). When we searched for brain structures whose responses are consistent
329 with the video's topic trajectory, we identified a network of structures that overlapped strongly
330 with the "long temporal receptive window" network reported by the Hasson group (e.g., compare
331 our Fig. 7A with the map of long temporal receptive window voxels in Lerner et al., 2011). This
332 provides support for the notion that part of the long temporal receptive window network may be
333 maintaining an explicit model of the stimulus dynamics. When we performed a similar analysis
334 after swapping out the video's topic trajectory with the recall topic trajectories of each individual
335 participant, this allowed us to identify brain regions whose responses (as the participants viewed
336 the video) reflected how the video trajectory would be transformed in memory (as reflected by
337 the recall topic trajectories). The analysis revealed that the rMTL and vmPFC may play a role in
338 this person-specific transformation from experience into memory. The role of the MTL in episodic
339 memory encoding has been well-reported (e.g., Paller and Wagner, 2002; Davachi et al., 2003;
340 Ranganath et al., 2004; Davachi, 2006; Wiltgen and Silva, 2007; Diana et al., 2007; van Kesteren
341 et al., 2013). Prior work has also implicated the medial prefrontal cortex in representing "schema"
342 knowledge (i.e., general knowledge about the format of an ongoing experience given prior similar
343 experiences; van Kesteren et al., 2012, 2013; Schlichting and Preston, 2015; Gilboa and Marlatte,
344 2017; Spalding et al., 2018). Integrating across our study and this prior work, one interpretation is
345 that the person-specific transformations mediated (or represented) by the rMTL and vmPFC may
346 reflect schema knowledge being leveraged, formed, or updated, incorporating ongoing experience
347 into previously acquired knowledge.

348 In extending classical free recall analyses to our naturalistic memory framework, we recovered
349 two patterns of recall dynamics central to list-learning studies: a high probability of initiating recall
350 with the first video event (Fig. 3A) and a strong bias toward transitioning from recalling a given

351 event to recalling the event immediately following it (Fig. 3B). However, equally noteworthy are
352 the typical free recall results not recovered in these analyses, as each highlights a fundamental
353 difference between list-learning studies and naturalistic memory paradigms like that employed
354 in the present study. The most noticeable departure from hallmark free recall dynamics in these
355 findings is the apparent lack of a serial position effect in Figure 3C, which instead shows greater and
356 lesser recall probabilities for events distributed across the video. Free recall experiments' stimuli
357 most often comprise lists of simple, common words, presented to participants in a random order.
358 (In fact, numerous word pools have been developed based on these criteria; e.g., Friendly et al.,
359 1982). These stimulus qualities enable two assumptions that are central to word list analyses, but
360 frequently do not hold for real-world experiences. First, researchers conducting free recall studies
361 may assume that the content at each presentation index is essentially equal, and does not bear
362 qualities that would affect participants' abilities to remember it above that of others. Such is rarely
363 the case with real-world experiences or experiments meant to approximate them, and the effects
364 of both intrinsic and observer-dependent factors on stimulus memorability are well-established
365 (for review see Chun and Turk-Browne, 2007; Bylinskii et al., 2015; Tyng et al., 2017). Second, the
366 random ordering of list items ensures that (across participants, on average) there is no relationship
367 between the thematic similarity of stimuli and their presentation positions—in other words, two
368 semantically related words are no more likely to be presented next to each other than at opposite
369 ends of the list. In most cases, the exact opposite is true of real-world episodes. Our internal
370 thoughts, our actions, and the physical state of the world around us all tend to follow a direct,
371 causal progression. As a result, each moment of our experience tends to be inherently more similar
372 to surrounding moments than to those in the distant past or future. Memory literature has termed
373 this strong temporal autocorrelation “context,” and in various media that depict real-world events
374 (e.g., movies and written stories), we recognize it as a *narrative structure*. While a random word
375 list (by definition) has no such structure, the logical progression between ideas and actions in
376 a naturalistic stimulus prompts the rememberer to recount presented events in order, from the
377 beginning. This tendency is reflected in our findings' second departure from typical free recall
378 dynamics: a lack of increased probability of first recall for end-of-sequence events (Fig. 3A).

379 Thus, analyses such as those in Figure 3 that address only the temporal dynamics of free re-
380 call paint an incomplete picture of memory for naturalistic episodes. While useful for studying
381 presentation order-dependent recall dynamics, they neglect to consider the stimuli’s content (or,
382 for example, that content’s potential interrelatedness). However, sensitivity to stimulus and recall
383 content introduces a new challenge: distinguishing between levels of recall quality for a stimulus
384 (i.e., an event) that is considered to have been “remembered.” When modeling memory experi-
385 ments, often times events (or items) and their later memories are treated as binary and independent
386 events (e.g., a given list item was simply either remembered or not remembered). Various models
387 of memory (e.g., Yonelinas, 2002) attempt to improve upon this by including confidence ratings,
388 rendering the binary judgement instead categorical. Our novel framework allows one to assess
389 memory performance in a more continuous way (*precision*), as well as analyze the correlational
390 structure of each encoding event to each memory event (*distinctiveness*). Further and importantly,
391 these two novel metrics we introduce here arise from comparisons of the actual content of the
392 experience/memories, which is not typically modeled. Leveraging this, we find that the successful
393 memory performance is related to 1) the precision with which the participant recounts each event
394 and 2) the distinctiveness of each recall event (relative to the other recalled events). The first finding
395 suggests that the information retained for *any individual event* may predict the overall amount of
396 information retained by the participant. The second finding suggests that the ability to distin-
397 guish between temporally or semantically similar content is also related to the quantity of content
398 recovered. Intriguingly, prior studies show that pattern separation, or the ability to discriminate
399 between similar experiences, is impaired in many cognitive disorders as well as natural aging
400 (Stark et al., 2010; Yassa et al., 2011; Yassa and Stark, 2011). Future work might explore whether
401 and how these metrics compare between cognitively impoverished groups and healthy controls.

402 While a large number of language models exist (e.g., WAS, LSA, word2vec, universal sentence
403 encoder; Steyvers et al., 2004; Landauer et al., 1998; Mikolov et al., 2013; Cer et al., 2018), here
404 we use latent dirichlet allocation (LDA)-based topic models for a few reasons. First, topic models
405 capture the *essence* of a text passage devoid of the specific set and order of words used. This was
406 an important feature of our model since different people may accurately recall a scene using very

407 different language. Second, words can mean different things in different contexts (e.g. “bat” as
408 the act of hitting a baseball, the object used for that action, or a flying mammal). Topic models
409 are robust to this, allowing words to exist as part of multiple topics. Last, topic models provide
410 a straightforward means to recover the weights for the particular words comprising a topic,
411 enabling easy interpretation of an event’s contents (e.g. Fig. 6). Other models such as Google’s
412 universal sentence encoder offer a context-sensitive encoding of text passages, but the encoding
413 space is complex and non-linear, and thus recovering the original words used to fit the model is
414 not straightforward. However, it’s worth pointing out that our framework is divorced from the
415 particular choice of language model. Moreover, many of the aspects of our framework could be
416 swapped out for other choices. For example, the language model, the timeseries segmentation
417 model and the video-recall matching function could all be customized for the particular problem.
418 Indeed for some problems, recovery of the particular recall words may not be necessary, and thus
419 other text-modeling approaches (such as universal sentence encoder) may be preferable. Future
420 work will explore the influence of particular model choices on the framework’s accuracy.

421 Our work has broad implications for how we characterize and assess memory in real-world
422 settings, such as the classroom or physician’s office. For example, the most commonly used
423 classroom evaluation tools involve simply computing the proportion of correctly answered exam
424 questions. Our work indicates that this approach is only loosely related to what educators might
425 really want to measure: how well did the students understand the key ideas presented in the
426 course? Under this typical framework of assessment, the same exam score of 50% could be
427 ascribed to two very different students: one who attended the full course but struggled to learn
428 more than a broad overview of the material, and one who attended only half of the course but
429 understood the material perfectly. Instead, one could apply our computational framework to build
430 explicit content models of the course material and exam questions. This approach would provide
431 a more nuanced and specific view into which aspects of the material students had learned well
432 (or poorly). In clinical settings, memory measures that incorporate such explicit content models
433 might also provide more direct evaluations of patients’ memories.

434 **Methods**

435 **Experimental design and data collection**

436 Data were collected by Chen et al. (2017). In brief, participants ($n = 17$) viewed the first 48 minutes
437 of “A Study in Pink”, the first episode of the BBC television series *Sherlock*, while fMRI volumes
438 were collected (TR = 1500 ms). The stimulus was divided into a 23 min (946 TR) and a 25 min
439 (1030 TR) segment to mitigate technical issues related to the scanner. After finishing the clip,
440 participants were instructed to (quoting from Chen et al., 2017) “describe what they recalled of the
441 [episode] in as much detail as they could, to try to recount events in the original order they were
442 viewed in, and to speak for at least 10 minutes if possible but that longer was better. They were told
443 that completeness and detail were more important than temporal order, and that if at any point
444 they realized they had missed something, to return to it. Participants were then allowed to speak
445 for as long as they wished, and verbally indicated when they were finished (e.g., ‘I’m done’).”
446 For additional details about the experimental procedure and scanning parameters, see Chen et al.
447 (2017). The experimental protocol was approved by Princeton University’s Institutional Review
448 Board.

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

455 **Data and code availability**

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

458 **Statistics**

459 All statistical tests we performed were two-sided.

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

461 **Topic modeling**

462 The input to the topic model we trained to characterize the dynamic content of the video comprised
463 hand-generated annotations of each of 1000 scenes spanning the video clip (generated by Chen
464 et al., 2017). The features annotated included: narrative details (a sentence or two describing
465 what happened in that scene); whether the scene took place indoors or outdoors; names of any
466 characters that appeared in the scene; name(s) of characters in camera focus; name(s) of characters
467 who were speaking in the scene; the location (in the story) that the scene took place; camera angle
468 (close up, medium, long, top, tracking, over the shoulder, etc.); whether music was playing in
469 the scene or not; and a transcription of any on-screen text. We concatenated the text for all of
470 these features within each segment, creating a “bag of words” describing each scene. We then
471 re-organized the text descriptions into overlapping sliding windows spanning 50 scenes each.
472 In other words, the first text sample comprised the combined text from the first 50 scenes (i.e.,
473 1–50), the second comprised the text from scenes 2–51, and so on. We trained our model using
474 these overlapping text samples with `scikit-learn` (version 0.19.1; Pedregosa et al., 2011), called
475 from our high-dimensional visualization and text analysis software, `HyperTools` (Heusser et al.,
476 2018b). Specifically, we used the `CountVectorizer` class to transform the text from each scene
477 into a vector of word counts (using the union of all words across all scenes as the “vocabulary,”
478 excluding English stop words); this yielded a number-of-scenes by number-of-words *word count*
479 matrix. We then used the `LatentDirichletAllocation` class (`topics=100, method='batch'`) to fit
480 a topic model (Blei et al., 2003) to the word count matrix, yielding a number-of-scenes (1000) by
481 number-of-topics (100) *topic proportions* matrix. The topic proportions matrix describes which mix
482 of topics (latent themes) is present in each scene. Next, we transformed the topic proportions
483 matrix to match the 1976 fMRI volume acquisition times. For each fMRI volume, we took the topic

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

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

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

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

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

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

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

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

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

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

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

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

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

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

537 **Probability of first recall (PFR).** PFR curves (Welch and Burnett, 1924; Postman and Phillips,
538 1965; Atkinson and Shiffrin, 1968) reflect the probability that an item will be recalled first as a
539 function of its serial position during encoding. To carry out this analysis, we initialized a number-
540 of-participants (17) by number-of-video-events (34) matrix of zeros. Then for each participant, we
541 found the index of the video event that was recalled first (i.e., the video event whose topic vector
542 was most strongly correlated with that of the first recall event) and filled in that index in the matrix
543 with a 1. Finally, we averaged over the rows of the matrix, resulting in a 1 by 34 array representing
544 the proportion of participants that recalled an event first, as a function of the order of the event's
545 appearance in the video (Fig. 3A).

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

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

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

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

575 **Semantic clustering scores.** Semantic clustering measures the extent to which participants clus-
576 tered their recall responses according to semantic similarity (Polyn et al., 2009). Here, we used the
577 topic vectors for each event as a proxy for its semantic content. Thus, the similarity between the
578 semantic content for two events can be computed by correlating their respective topic vectors. For
579 each recall event transition, we sorted all not-yet-recalled events according to how correlated the
580 topic vector of *the closest-matching video event* was to the topic vector of the closest-matching video
581 event to the just-recalled event. We then computed the percentile rank of the observed next recall.
582 We averaged these percentile ranks across all of the participant's recalls to obtain a single semantic
583 clustering score for the participant (mean: 0.813, SEM: 0.022). We found that participants who
584 exhibited stronger semantic clustering scores overall remembered more video events (Pearson's
585 $r(15) = 0.55, p = 0.02$).

586 **Novel naturalistic memory metrics**

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

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

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

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

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

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

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

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

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

634 **Searchlight fMRI analyses**

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

642 Next, we constructed two sets of “template” matrices: one reflected the video’s topic trajectory
643 and the other reflected each participant’s recall topic trajectory. To construct the video template, we
644 computed the correlations between the topic proportions estimated for every pair of TRs (prior to
645 segmenting the trajectory into discrete events; i.e., the correlation matrix shown in Figs. 2B and 7A).
646 We constructed similar temporal correlation matrices for each participant’s recall topic trajectory
647 (Figs. 2D, S4). However, to correct for length differences and potential non-linear transformations
648 between viewing time and recall time, we first used dynamic time warping (Berndt and Clifford,
649 1994) to temporally align participants’ recall topic trajectories with the video topic trajectory (an
650 example correlation matrix before and after warping is shown in Fig. 7B). This yielded a 1976 by
651 1976 correlation matrix for the video template and for each participant’s recall template.

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

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

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

672 References

- 673 Atkinson, R. C. and Shiffrin, R. M. (1968). Human memory: A proposed system and its control
674 processes. In Spence, K. W. and Spence, J. T., editors, *The psychology of learning and motivation*,
675 volume 2, pages 89–105. Academic Press, New York.
- 676 Baldassano, C., Chen, J., Zadbood, A., Pillow, J. W., Hasson, U., and Norman, K. A. (2017).
677 Discovering event structure in continuous narrative perception and memory. *Neuron*, 95(3):709–
678 721.
- 679 Berndt, D. J. and Clifford, J. (1994). Using dynamic time warping to find patterns in time series. In
680 *KDD workshop*, volume 10, pages 359–370.
- 681 Blei, D. M. and Lafferty, J. D. (2006). Dynamic topic models. In *Proceedings of the 23rd International
682 Conference on Machine Learning*, ICML '06, pages 113–120, New York, NY, US. ACM.
- 683 Blei, D. M., Ng, A. Y., and Jordan, M. I. (2003). Latent dirichlet allocation. *Journal of Machine
684 Learning Research*, 3:993 – 1022.

- 685 Brunec, I. K., Moscovitch, M. M., and Barense, M. D. (2018). Boundaries shape cognitive representations of spaces and events. *Trends in Cognitive Sciences*, 22(7):637–650.
- 686
- 687 Bylinskii, Z., Isola, P., Bainbridge, C., Torralba, A., and Oliva, A. (2015). Intrinsic and extrinsic
688 effects on image memorability. *Vision Research*, 116:165–178.
- 689 Capota, M., Turek, J., Chen, P.-H., Zhu, X., Manning, J. R., Sundaram, N., Keller, B., Wang, Y., and
690 Shin, Y. S. (2017). Brain imaging analysis kit.
- 691 Cer, D., Yang, Y., y Kong, S., Hua, N., Limtiaco, N., John, R. S., Constant, N., Guajardo-Cespedes,
692 M., Yuan, S., Tar, C., Sung, Y.-H., Strope, B., and Kurzweil, R. (2018). Universal sentence encoder.
693 *arXiv*, 1803.11175.
- 694 Chen, J., Leong, Y. C., Honey, C. J., Yong, C. H., Norman, K. A., and Hasson, U. (2017). Shared
695 memories reveal shared structure in neural activity across individuals. *Nature Neuroscience*,
696 20(1):115.
- 697 Chun, M. and Turk-Browne, N. (2007). Interactions between attention and memory. *Current opinion
698 in neurobiology*, 17(2):177–184.
- 699 Clewett, D. and Davachi, L. (2017). The ebb and flow of experience determines the temporal
700 structure of memory. *Curr Opin Behav Sci*, 17:186–193.
- 701 Davachi, L. (2006). Item, context and relational episodic encoding in humans. *Current Opinion in
702 Neurobiology*, 16(6):693—700.
- 703 Davachi, L., Mitchell, J. P., and Wagner, A. D. (2003). Multiple routes to memory: distinct medial
704 temporal lobe processes build item and source memories. *Proceedings of the National Academy of
705 Sciences, USA*, 100(4):2157 – 2162.
- 706 Diana, R. A., Yonelinas, A. P., and Ranganath, C. (2007). Imaging recollection and famil-
707 iarity in the medial temporal lobe: a three-component model. *Trends in Cognitive Sciences*,
708 doi:10.1016/j.tics.2007.08.001.

- 709 DuBrow, S. and Davachi, L. (2013). The influence of contextual boundaries on memory for the
710 sequential order of events. *Journal of Experimental Psychology: General*, 142(4):1277–1286.
- 711 Ezzyat, Y. and Davachi, L. (2011). What constitutes an episode in episodic memory? *Psychological
712 Science*, 22(2):243–252.
- 713 Friendly, M., Franklin, P. E., Hoffman, D., and Rubin, D. C. (1982). The Toronto Word Pool:
714 Norms for imagery, concreteness, orthographic variables, and grammatical usage for 1,080
715 words. *Behavior Research Methods and Instrumentation*, 14:375–399.
- 716 Gilboa, A. and Marlatte, H. (2017). Neurobiology of schemas and schema-mediated memory.
717 *Trends Cogn Sci*, 21(8):618–631.
- 718 Heusser, A. C., Ezzyat, Y., Shiff, I., and Davachi, L. (2018a). Perceptual boundaries cause mnemonic
719 trade-offs between local boundary processing and across-trial associative binding. *Journal of
720 Experimental Psychology Learning, Memory, and Cognition*, 44(7):1075–1090.
- 721 Heusser, A. C., Fitzpatrick, P. C., Field, C. E., Ziman, K., and Manning, J. R. (2017). Quail: a
722 Python toolbox for analyzing and plotting free recall data. *The Journal of Open Source Software*,
723 10.21105/joss.00424.
- 724 Heusser, A. C., Ziman, K., Owen, L. L. W., and Manning, J. R. (2018b). HyperTools: a Python
725 toolbox for gaining geometric insights into high-dimensional data. *Journal of Machine Learning
726 Research*, 18(152):1–6.
- 727 Howard, M. W. and Kahana, M. J. (2002). A distributed representation of temporal context. *Journal
728 of Mathematical Psychology*, 46:269–299.
- 729 Howard, M. W., MacDonald, C. J., Tiganj, Z., Shankar, K. H., Du, Q., Hasselmo, M. E., and H., E.
730 (2014). A unified mathematical framework for coding time, space, and sequences in the medial
731 temporal lobe. *Journal of Neuroscience*, 34(13):4692–4707.
- 732 Howard, M. W., Viskontas, I. V., Shankar, K. H., and Fried, I. (2012). Ensembles of human MTL
733 neurons “jump back in time” in response to a repeated stimulus. *Hippocampus*, 22:1833–1847.

- 734 Huk, A., Bonnen, K., and He, B. J. (2018). Beyond trial-based paradigms: continuous behavior, on-
735 going neural activity, and naturalistic stimuli. *Journal of Neuroscience*, 10.1523/JNEUROSCI.1920-
736 17.2018.
- 737 Kahana, M. J. (1996). Associative retrieval processes in free recall. *Memory & Cognition*, 24:103–109.
- 738 Koriat, A. and Goldsmith, M. (1994). Memory in naturalistic and laboratory contexts: distin-
739 guishing accuracy-oriented and quantity-oriented approaches to memory assessment. *Journal of*
740 *Experimental Psychology: General*, 123(3):297–315.
- 741 Kriegeskorte, N., Mur, M., and Bandettini, P. (2008). Representational similarity analysis – con-
742 necting the branches of systems neuroscience. *Frontiers in Systems Neuroscience*, 2:1 – 28.
- 743 Landauer, T. K., Foltz, P. W., and Laham, D. (1998). Introduction to latent semantic analysis.
744 *Discourse Processes*, 25:259–284.
- 745 Lerner, Y., Honey, C. J., Silbert, L. J., and Hasson, U. (2011). Topographic mapping of a hierarchy
746 of temporal receptive windows using a narrated story. *Journal of Neuroscience*, 31(8):2906–2915.
- 747 Manning, J. (2019). Episodic memory: mental time travel or a quantum ‘memory wave’ function?
748 *PsyArXiv*, doi:10.31234/osf.io/6zjwb.
- 749 Manning, J. R., Norman, K. A., and Kahana, M. J. (2015). The role of context in episodic memory.
750 In Gazzaniga, M., editor, *The Cognitive Neurosciences, Fifth edition*, pages 557–566. MIT Press.
- 751 Manning, J. R., Polyn, S. M., Baltuch, G., Litt, B., and Kahana, M. J. (2011). Oscillatory patterns
752 in temporal lobe reveal context reinstatement during memory search. *Proceedings of the National*
753 *Academy of Sciences, USA*, 108(31):12893–12897.
- 754 McInnes, L. and Healy, J. (2018). UMAP: Uniform manifold approximation and projection for
755 dimension reduction. *arXiv*, 1802(03426).
- 756 Mikolov, T., Chen, K., Corrado, G., and Dean, J. (2013). Efficient estimation of word representations
757 in vector space. *arXiv*, 1301.3781.

- 758 Mueller, A., Fillion-Robin, J.-C., Boidol, R., Tian, F., Nechifor, P., yoonsubKim, Peter, Rampin, R.,
759 Corvellec, M., Medina, J., Dai, Y., Petrushev, B., Langner, K. M., Hong, Alessio, Ozsvald, I.,
760 vkolmakov, Jones, T., Bailey, E., Rho, V., IgorAPM, Roy, D., May, C., foobuzz, Piyush, Seong,
761 L. K., Goey, J. V., Smith, J. S., Gus, and Mai, F. (2018). WordCloud 1.5.0: a little word cloud
762 generator in Python. *Zenodo*, <https://zenodo.org/record/1322068#.W4tPKZNKh24>.
- 763 Murdock, B. B. (1962). The serial position effect of free recall. *Journal of Experimental Psychology*,
764 64:482–488.
- 765 Paller, K. A. and Wagner, A. D. (2002). Observing the transformation of experience into memory.
766 *Trends in Cognitive Sciences*, 6(2):93–102.
- 767 Pedregosa, F., Varoquaux, G., Gramfort, A., Michel, V., Thirion, B., Grisel, O., Blondel, M., Pretten-
768 hofer, P., Weiss, R., Dubourg, V., Vanderplas, J., Passos, A., Cournapeau, D., Brucher, M., Perrot,
769 M., and Duchesnay, E. (2011). Scikit-learn: Machine learning in Python. *Journal of Machine
770 Learning Research*, 12:2825–2830.
- 771 Polyn, S. M., Norman, K. A., and Kahana, M. J. (2009). A context maintenance and retrieval model
772 of organizational processes in free recall. *Psychological Review*, 116(1):129–156.
- 773 Postman, L. and Phillips, L. W. (1965). Short-term temporal changes in free recall. *Quarterly Journal
774 of Experimental Psychology*, 17:132–138.
- 775 Rabiner, L. (1989). A tutorial on Hidden Markov Models and selected applications in speech
776 recognition. *Proceedings of the IEEE*, 77(2):257–286.
- 777 Radvansky, G. A. and Zacks, J. M. (2017). Event boundaries in memory and cognition. *Curr Opin
778 Behav Sci*, 17:133–140.
- 779 Ranganath, C., Cohen, M. X., Dam, C., and D’Esposito, M. (2004). Inferior temporal, prefrontal,
780 and hippocampal contributions to visual working memory maintenance and associative memory
781 retrieval. *Journal of Neuroscience*, 24(16):3917–3925.

- 782 Ranganath, C. and Ritchey, M. (2012). Two cortical systems for memory-guided behavior. *Nature*
783 *Reviews Neuroscience*, 13:713 – 726.
- 784 Schlichting, M. L. and Preston, A. R. (2015). Memory integration: neural mechanisms and impli-
785 cations for behavior. *Current Opinion in Behavioral Sciences*, 1:1–8.
- 786 Simony, E., Honey, C. J., Chen, J., and Hasson, U. (2016). Uncovering stimulus-locked network
787 dynamics during narrative comprehension. *Nature Communications*, 7(12141):1–13.
- 788 Spalding, K. N., Schlichting, M. L., Zeithamova, D., Preston, A. R., Tranel, D., Duff, M. C., and
789 Warren, D. E. (2018). Ventromedial prefrontal cortex is necessary for normal associative inference
790 and memory integration. *The Journal of Neuroscience*, 38(15):3767–3775.
- 791 Stark, S. M., Yassa, M. A., and Stark, C. E. L. (2010). Individual differences in spatial pattern
792 separation performance associated with healthy aging in humans. *Learning & Memory*, 17(6):284–
793 288.
- 794 Steyvers, M., Shiffrin, R. M., and Nelson, D. L. (2004). Word association spaces for predicting
795 semantic similarity effects in episodic memory. In Healy, A. F., editor, *Cognitive Psychology and*
796 *its Applications: Festschrift in Honor of Lyle Bourne, Walter Kintsch, and Thomas Landauer*. American
797 Psychological Association, Washington, DC.
- 798 Tompry, A. and Davachi, L. (2017). Consolidation promotes the emergence of representational
799 overlap in the hippocampus and medial prefrontal cortex. *Neuron*, 96(1):228–241.
- 800 Tyng, C. M., Amin, H. U., Saad, M. N. M., and S, M. A. (2017). The influences of emotion on
801 learning and memory. *Frontiers in psychology*, 8:1454.
- 802 van Kesteren, M. T. R., Beul, S. F., Takashima, A., Henson, R. N., Ruiter, D. J., and Fernández, G.
803 (2013). Differential roles for medial prefrontal and medial temporal cortices in schema-dependent
804 encoding: from congruent to incongruent. *Neuropsychologia*, 51(12):2352–2359.
- 805 van Kesteren, M. T. R., Ruiter, D. J., Fernández, G., and Henson, R. N. (2012). How schema and
806 novelty augment memory formation. *Trends Neurosci*, 35(4):211–9.

- 807 Waskom, M., Botvinnik, O., Okane, D., Hobson, P., David, Halchenko, Y., Lukauskas, S., Cole, J. B.,
808 Warmenhoven, J., de Ruiter, J., Hoyer, S., Vanderplas, J., Villalba, S., Kunter, G., Quintero, E.,
809 Martin, M., Miles, A., Meyer, K., Augspurger, T., Yarkoni, T., Bachant, P., Williams, M., Evans,
810 C., Fitzgerald, C., Brian, Wehner, D., Hitz, G., Ziegler, E., Qalieh, A., and Lee, A. (2016). Seaborn:
811 v0.7.1.
- 812 Welch, G. B. and Burnett, C. T. (1924). Is primacy a factor in association-formation. *American Journal
813 of Psychology*, 35:396–401.
- 814 Wiltgen, B. J. and Silva, A. J. (2007). Memory for context becomes less specific with time. *Learning
815 & Memory*, 14(4):313–317.
- 816 Yassa, M. A., Lacy, J. W., Stark, S. M., Albert, M. S., Gallagher, M., and Stark, C. E. L. (2011). Pattern
817 separation deficits associated with increased hippocampal ca3 and dentate gyrus activity in
818 nondemented older adults. *Hippocampus*, 21(9):968–979.
- 819 Yassa, M. A. and Stark, C. E. L. (2011). Pattern separation in the hippocampus. *Trends In Neuro-
820 sciences*, 34(10):515–525.
- 821 Yonelinas, A. P. (2002). The nature of recollection and familiarity: A review of 30 years of research.
822 *Journal of Memory and Language*, 46:441–517.
- 823 Zacks, J. M., Speer, N. K., Swallow, K. M., Braver, T. S., and Reynolds, J. R. (2007). Event perception:
824 a mind-brain perspective. *Psychological Bulletin*, 133:273–293.
- 825 Zadbood, A., Chen, J., Leong, Y. C., Norman, K. A., and Hasson, U. (2017). How we transmit
826 memories to other brains: Constructing shared neural representations via communication. *Cereb
827 Cortex*, 27(10):4988–5000.
- 828 Zwaan, R. A. and Radvansky, G. A. (1998). Situation models in language comprehension and
829 memory. *Psychological Bulletin*, 123(2):162 – 185.

830 **Supporting information**

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

832 **Acknowledgements**

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

838 **Author contributions**

839 Conceptualization: A.C.H. and J.R.M.; Methodology: A.C.H. and J.R.M.; Software: A.C.H., P.C.F.
840 and J.R.M.; Analysis: A.C.H., P.C.F. and J.R.M.; Writing, Reviewing, and Editing: A.C.H., P.C.F.
841 and J.R.M.; Supervision: J.R.M.

842 **Author information**

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