# *Supporting Information for*: How is experience transformed into memory?

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

September 5, 2018

## Overview

This document provides additional details about the methods we used in the main text. We also include some additional analyses and figures referenced in the main text.

## Additional details about topic modeling methods and results

### Optimizing topic model parameters

In order to create accurate video and recall models, we used an optimization method that was driven by our ability to explain hand-annotated memory performance metrics collected by Chen et al. (2017). Specifically, we used a grid search to compute the $\omega$ (video sliding window duration, in scenes), $\rho$ (recall sliding window duration, in sentences), and $K$ (number of topics) that satisfied

$$\underset{\omega,\rho,K}{\mathrm{argmax}} \left[ \mathrm{corr} \left( \mathrm{corr} \left( \mu \left( \omega, \rho, K \right), \nu \left( \omega, \rho, K \right) \right), \theta \right) \right],$$

where $\mathrm{corr}(\mu, \nu)$ is the per-participant correlation between the upper triangles of the temporal correlation matrices of the video ($\mu$) and recall ($\nu$) trajectory, and $\theta$ is the per-participant hand-annotated memory performance. We searched over a grid of pre-specified values for each of these parameters; the resulting correlations are displayed in Figure S1. The optimal parameters were $\omega = 50$, $\rho = 10$, and $K = 100$.

The optimized model converged on 27 unique topics that were assigned non-zero weights over the course of the video. We provide a list of the top ten highest-weighted words from each topic in Figure S2.

### Feature importance analyses

To determine the contribution of each feature to the structure of the video topic proportions, we conducted a "leave one out" analysis. Specifically, we compared the original video topic trajectory
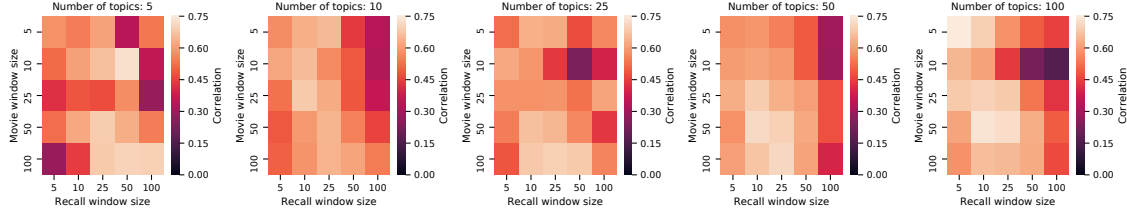
1

**Figure S1: Optimizing topic model parameters.** We performed a grid search over video sliding window length ($\omega \in \{5, 10, 25, 50, 100\}$), recall sliding window length ($\rho \in \{5, 10, 25, 50, 100\}$, and number of topics ($K \in \{5, 10, 25, 50, 100\}$. The reported correlations are between per-subject video-recall trajectory correlations and per-subject hand-annotated memory performance ratings.

(created using all hand-annotated features from the 1000 hand-annotated scenes spanning the *Sherlock* episode; see *Methods* for a full list of features) with video trajectories created using all but one type of feature. We created temporal correlation matrices for each trajectory (using the topic proportions matrices) and correlated the upper triangles of each impoverished trajectory with the original feature-complete trajectory. Observing a lower correlation between an impoverished trajectory (holding out a particular feature) and the feature-complete trajectory would suggest that the given feature played a more prominent role in shaping the structure of the feature-complete trajectory. We found that hand-annotated narrative details provided the most structure to the feature-complete trajectory, whereas transcriptions of onscreen text provided the least structure (Fig. S3A).

We also carried out an analysis of which annotated features tended to shape aspects of the video topic trajectory that were preserved in participants' recalls. Specifically, we computed the timepoint-by-timepoint correlation matrix of the video topic trajectory, and correlated its upper triangle with that of the timepoint-by-timepoint correlation matrices of each participant's recall topic trajectory (resampled using linear interpolation to have the same number of timepoints as the video trajectory). This yielded a single correlation coefficient for each participant. We then repeated this analysis with each annotated feature held out in turn. Observing a lower correlation between the video and recall trajectories (when a given feature was held out) would indicate that the feature tends to be preserved in participants' recalls. We found that hand-annotated narrative details were the most preserved type of feature, whereas information about the camera angle tended not to influence participants' recalls (Fig. S3B).

Next, we wondered how the different types of features might relate. For example, knowing which characters are on screen during a given scene may also provide information about which characters are speaking. We computed video topic trajectories for each feature in turn, and then compared the temporal correlation matrices of all pairs of features. This provided additional confirmation that the shape of the full trajectory (including all types of features) was largely driven by narrative details. We also found that character-driven features (characters on screen, characters speaking, and characters in focus) were strongly correlated. Other details, such as the presence or absence of music, led to very different topic trajectories (Fig. S3C).

| Topic ID | Top 10 words | Topic description |
|---|---|---|
| 1 | sir, jeffrey, indoor, yes, office, building, aide, helen, lestrade, medium | The first death |
| 2 | sherlock, john, outdoor, taxi, yes, medium, road, says, phone, continues | John being followed (a) |
| 3 | sherlock, john, donovan, medium, lauriston, gardens, anderson, street, outdoor, lestrade | Discussing the fourth death |
| 4 | lestrade, donovan, room, indoor, press, conference, police, medium, reporter, reporters | Press conference (a) |
| 5 | john, man, yes, warehouse, indoor, medium, shoulder, says, hand, asks | Meeting with Mycroft (a) |
| 6 | sherlock, lestrade, john, indoor, medium, gardens, lauriston, room, floor, crime | Examining a body (a) |
| 7 | john, road, brixton, outdoor, phone, box, yes, medium, man, camera | John being followed (b) |
| 8 | john, sherlock, street, medium, baker, indoor, says, mrs, hudson, 221b | 221b Baker St. (a) |
| 9 | john, donovan, lauriston, gardens, yes, street, medium, outdoor, shoulder, policeman | Consulting with the police |
| 10 | lestrade, donovan, indoor, room, medium, aide, press, conference, police, reporter | Press conference (b) |
| 11 | john, mike, lestrade, medium, donovan, park, indoor, square, russell, outdoor | Exposition |
| 12 | john, sherlock, medium, street, baker, anthea, indoor, yes, 221b, suite | Bringing John back |
| 13 | sherlock, john, st, bartholomew, hospital, indoor, medium, molly, mike, laboratory | John meets Sherlock (a) |
| 14 | john, man, yes, anthea, medium, warehouse, indoor, car, road, outdoor | Kidnapping John |
| 15 | john, mike, sherlock, medium, molly, park, russell, square, outdoor, bench | John runs into an old friend |
| 16 | jimmy, yes, indoor, donovan, medium, aide, gary, lestrade, press, conference | The second death (a) |
| 17 | sherlock, john, crime, scene, room, floor, lauriston, gardens, indoor, lestrade | Examining a body (b) |
| 18 | sherlock, john, mrs, hudson, baker, street, 221b, indoor, suite, yes | 221b Baker St. (b) |
| 19 | john, jeffrey, sir, indoor, yes, medium, psychotherapist, helen, office, london | John's psychotherapy appointment |
| 20 | john, sherlock, yes, laboratory, indoor, hospital, bartholomew, st, medium, mike | John meets sherlock (b) |
| 21 | sherlock, lestrade, indoor, yes, room, floor, gardens, lauriston, scene, crime | Examining a body (c) |
| 22 | john, indoor, room, medium, psychotherapist, yes, soldiers, close, london, outdoor | John's PTSD |
| 23 | yes, jeffrey, sir, jimmy, aide, indoor, medium, woman, helen, man | Press conference (c) |
| 24 | sherlock, john, suite, street, 221b, baker, indoor, medium, says, asks | 221b Baker St. (c) |
| 25 | man, john, warehouse, indoor, yes, shoulder, medium, says, continues, looks | Meeting with Mycroft (b) |
| 26 | jimmy, yes, gary, sir, jeffrey, medium, indoor, outdoor, psychotherapist, rain | The second death (b) |
| 27 | sherlock, john, indoor, street, baker, medium, 221b, suite, yes, phone | 221b Baker St. (d) |

**Figure S2: Topics discovered in *Sherlock*.** We applied a topic model to hand-annotated information about 1000 scenes spanning the 45 minute episode. We identified 27 unique topics with non-zero weights (we used $K = 100$ topics to fit the model). Each topic comprises a distribution of weights over all words in the vocabulary. For each topic, we show the words with the 10 largest weights, along with a suggested description of the topic.
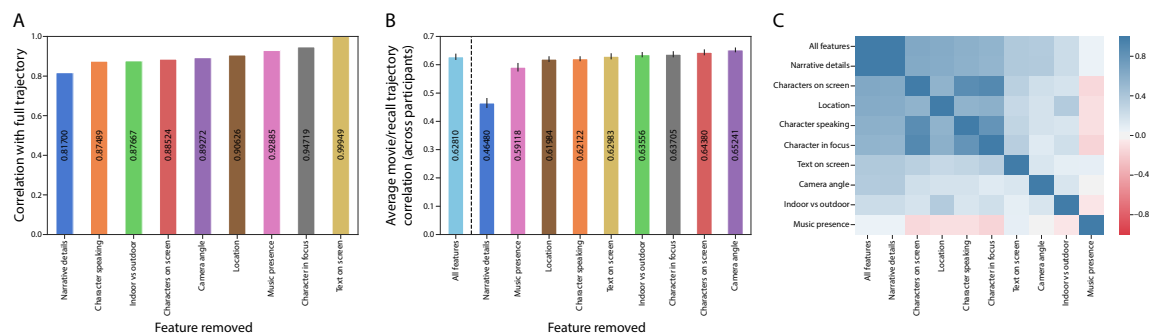
**Figure S3: Feature importance analysis. A.** Contributions of each feature type to the structure of the video trajectory. The bar heights reflect the correlation between the video trajectory computed using all features with a video trajectory computed using all features except the indicated feature. (Lower bars reflect features that contribute more substantially to the video trajectory's shape.) **B.** Which features are preserved during recall? The bar heights reflect the (average) across-participant correlations between the video and recall trajectories. Error bars denote bootstrap-estimated standard error of the mean. **C.** Feature correlation matrix. Each entry displays the correlation between video topic trajectories created using only the indicated (row/column) features.

# Additional analyses of memory performance

## Naturalistic extensions of classic list-learning analyses

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

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

As described below, we next considered three more nuanced measures of the memory performance and dynamics that are typically associated with list-learning studies. We also provide a software package, Quail, for carrying out these analyses (Heusser et al., 2017).

**Probability of first recall (PFR).** PFR curves (Welch and Burnett, 1924; Postman and Phillips, 1965; Atkinson and Shiffrin, 1968) 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 (34) 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
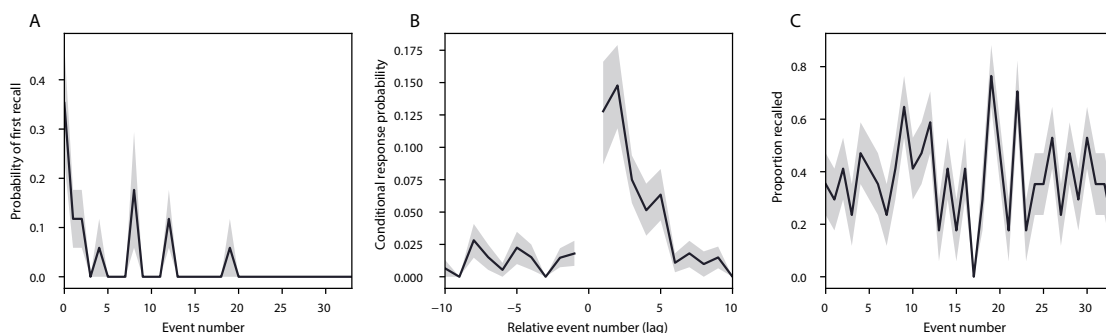
**Figure S4: Naturalistic extensions of classic list-learning memory analyses. A.** The probability of first recall as a function of the serial position of the event in the video. **B**. The probability of recalling each event, conditioned on having most recently recalled the event *lag* events away 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.

with a 1. Finally, we averaged over the rows of the matrix, resulting in a 1 by 34 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. S4A).

**Lag conditional probability curve (lag-CRP).** The lag-CRP curve (Kahana, 1996) reflects the probability of recalling a given event after the just-recalled event, as a function of their relative positions (or *lag*). In other words, a lag of 1 indicates that a recalled event came immediately after the previously recalled event in the video, and a lag of -3 indicates that a recalled event came 3 events before the previously recalled event. 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 (-33 to +33; 67 lags total) matrix. We averaged over the rows of this matrix to obtain a group-averaged lag-CRP curve (Fig. S4B).

**Serial position curve (SPC).** SPCs (Murdock, 1962) reflect the proportion of participants that remember each item as a function of their serial position during encoding. We initialized a number-of-participants (17) by number-of-video-events (34) 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 (i.e., for the given participant and event). 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 34 array representing the proportion of participants that recalled each event as a function of the order of the event's appearance in the video (Fig. S4C).

**Temporal clustering scores.** Temporal clustering refers to the extent to which participants group their recall responses according to encoding position (Polyn et al., 2009). For instance, if a participant recalled the video events in the exact order they occurred (or in exact reverse order), this

would yield a score of 1. If a participant recalled the events in random order, this would yield an expected score of 0.5. For each recall event transition (and separately for each participant), we sorted all not-yet-recalled events according to their absolute lag (i.e., distance away in the video). We then computed the percentile rank of the next event the participant recalled. We averaged these percentile ranks across all of the participant's recalls to obtain a single temporal clustering score for the participant (mean: 0.808, SEM: 0.022). Overall, we found that participants with higher temporal clustering scores also tended to recall more events ($r = 0.62, p = 0.007$).

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

## Additional measures of naturalistic memory

To quantify the similarity between the video topic trajectory and individual recall topic trajectories, we considered several novel metrics. First, we tested whether each participant's recall trajectory matched the video trajectory in a general sense. For each participant we filtered the video trajectory to only include the events that the participant remembered. We then computed the root mean squared difference (RMSD) between the remaining video events and the (closest-matching) recalled events. For example, if the topic vectors for a participant's recall event topic vectors matched the corresponding video event topic vectors exactly (and in order), the expected RMSD for those events would be 0. However, if the participant's recall events did not perfectly match the video events, or if they were out of order, then the RMSD would be greater than 0. To assess the significance of the match between the video and recall trajectories, we carried out a permutation procedure whereby, for each of 10000 repetitions, we circularly shifted the recall trajectories (in time) by a random amount and then re-computed the RMSD each time. This yielded a distribution of "null" RMSD values for each participant. The observed RMSD values reached significance (i.e., $p < 0.05$, reflecting that more than 95% of the null RMSD values were greater than the observed RMSD value) for nine of the participants (3, 4, 8–13, and 17). (For the remaining participants this test yielded $0.05 \leq p < 0.25$.) The observed RMSD values were also reliably correlated with hand-annotated memory performance across participants ($r = -0.57, p = 0.016$). In other words, a closer match between the video and recall topic trajectories was related to better overall recall performance.

**Precision.** We next tested whether participants who recalled more events were also more *precise* in their recollections. For each participant, we computed the correlation between the topic vectors for each recall event and that of its closest-matching video event (only for the events which they recalled). We defined the precision as the average video-recall correlation across all of the events a participant recalled. We found a strong correlation between hand-annotated memory performance

and precision, suggesting that participants who remembered more events also remembered them more veridically ($r = 0.74, p = 0.0006$).

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

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

where $\bar{c}(\text{event})$ is the average correlation between the given recalled event's topic vector and the topic vectors from all video events *except* the best-matching video event. We then averaged these distinctiveness scores across all of the events recalled by the given participant. We found that participants with higher average distinctiveness scores tended to also have better hand-annotated memory performance ($r = 0.8, p = 0.0001$).

**Other order effects.** We tested whether participants with better memory performance were also more likely to remember the events in order. For each participant, we computed the Spearman rank correlation between the order of events that the participant recalled and the order the events actually occurred in the video (considering in the analysis only the events that the participant recalled). Participants who recalled more events also recalled more of them in order ($r = 0.5, p = 0.04$). In summary, we found that better memory performance was associated with more precise, distinctive, and ordered recalls.

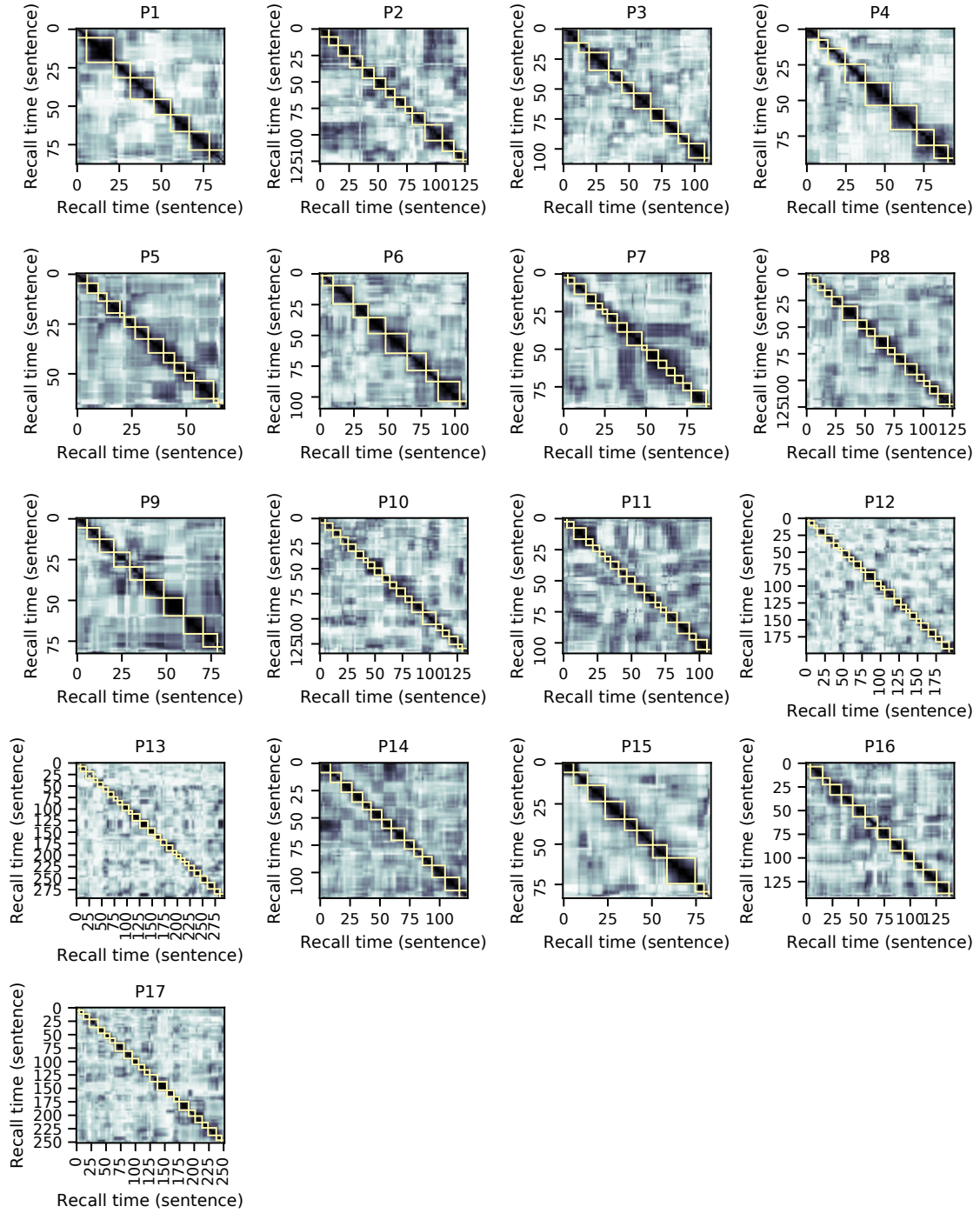# Participant-level figures referenced in the main text

**Figure S5: Recall trajectory temporal correlation matrices and event segmentation fits.** Each panel is in the same format as Figure 2E in the main text. The yellow boxes indicate HMM-identified event boundaries.
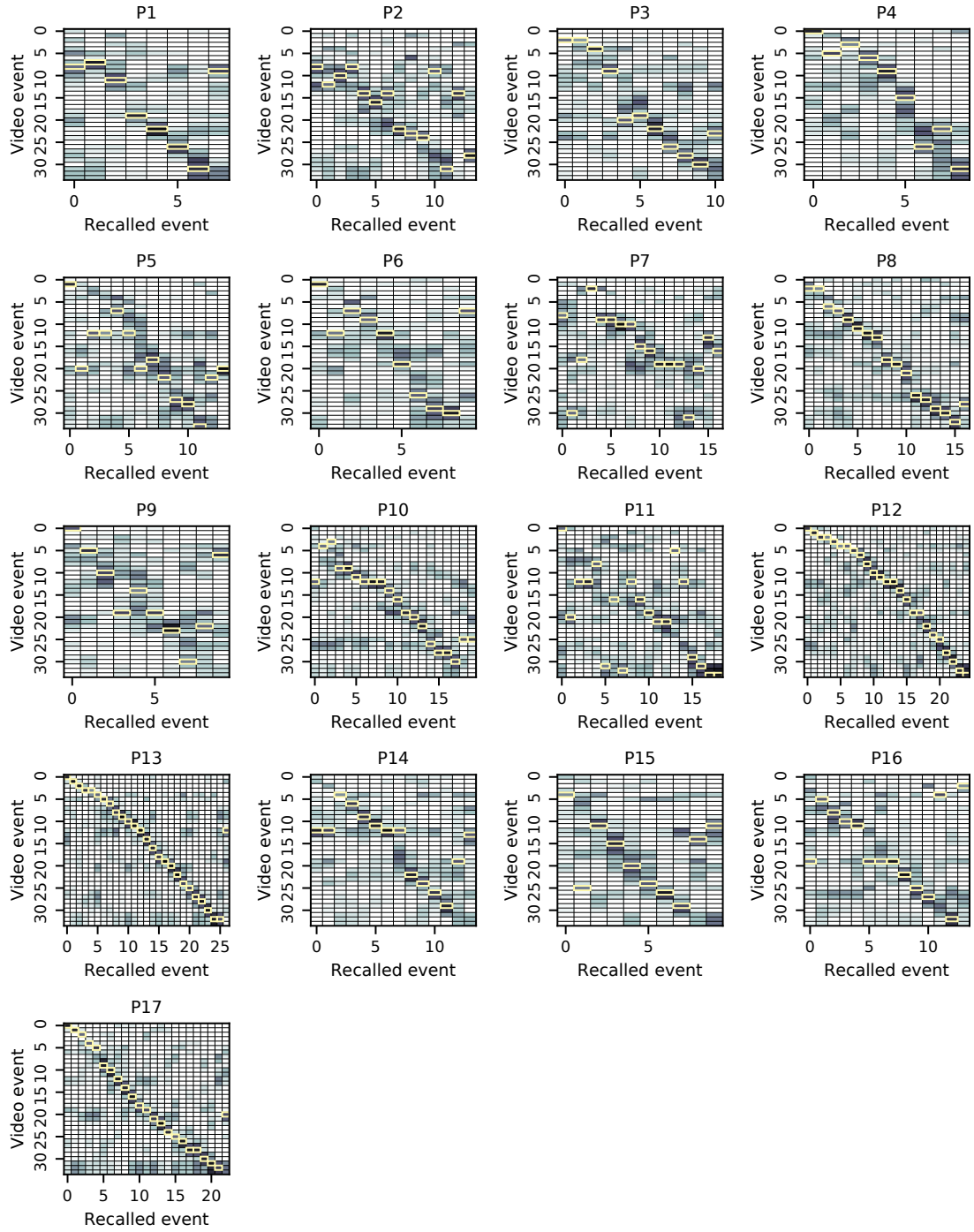
8

**Figure S6: Video-recall event correlation matrices.** Each panel is in the same format as Figure 2G in the main text. The yellow boxes mark the maximum correlation in each column.

# Supplemental references

Atkinson, R. C. and Shiffrin, R. M. (1968). Human memory: A proposed system and its control processes. In Spence, K. W. and Spence, J. T., editors, *The psychology of learning and motivation*, volume 2, pages 89–105. Academic Press, New York.

Chen, J., Leong, Y. C., Norman, K. A., and Hasson, U. (2017). Shared experience, shared memory: a common structure for brain activity during naturalistic recall shared experience, shared memory: a common structure for brain activity during naturalistic recall. *Nature Neuroscience*, 20(1):115–125.

Heusser, A. C., Fitzpatrick, P. C., Field, C. E., Ziman, K., and Manning, J. R. (2017). Quail: a Python toolbox for analyzing and plotting free recall data. *The Journal of Open Source Software*, 10.21105/joss.00424.

Kahana, M. J. (1996). Associative retrieval processes in free recall. *Memory & Cognition*, 24:103–109.

Murdock, B. B. (1962). The serial position effect of free recall. *Journal of Experimental Psychology*, 64:482–488.

Polyn, S. M., Norman, K. A., and Kahana, M. J. (2009). A context maintenance and retrieval model of organizational processes in free recall. *Psychological Review*, 116(1):129–156.

Postman, L. and Phillips, L. W. (1965). Short-term temporal changes in free recall. *Quarterly Journal of Experimental Psychology*, 17:132–138.

Welch, G. B. and Burnett, C. T. (1924). Is primacy a factor in association-formation. *American Journal of Psychology*, 35:396–401.