

Neural measures of subsequent memory reflect endogenous variability in cognitive function

Christoph T. Weidemann¹ & Michael J. Kahana²

¹*Swansea University*

²*University of Pennsylvania*

Humans cognition exhibits a striking degree of variability: Sometimes we rapidly forge new associations whereas at others new information simply does not stick. Although strong correlations between neural activity during encoding and subsequent retrieval performance have implicated such “subsequent memory effects” (SMEs) as important for understanding the neural basis of memory formation, uncontrolled variability in external factors that also predict memory performance confounds the interpretation of these effects. By controlling for a comprehensive set of external variables, we investigated the extent to which neural correlates of successful memory encoding reflect variability in endogenous brain states. We show that external variables that reliably predict memory performance have only minimal effects on electroencephalographic (EEG) correlates of successful memory encoding. Instead, the brain activity that is diagnostic of successful encoding primarily reflects fluctuations in endogenous neural activity. These findings link neural activity during learning to endogenous states that drive variability in human cognition.

The capacity to learn new information can vary considerably from moment to moment. We all recognize this variability in the frustration and embarrassment that accompanies associated

memory lapses. Researchers investigate the neural basis of this variability by analyzing brain activity during the encoding phase of a memory experiment as a function of each item's subsequent retrieval success. Across hundreds of such studies, the resulting contrasts, termed subsequent memory effects (SMEs), have revealed reliable biomarkers of successful memory encoding.¹⁻³

A key question, however, is whether the observed SMEs actually indicate endogenously varying brain states, or whether they instead reflect variation in external stimulus- and task-related variables, such as item difficulty or proactive interference, known to strongly predict retrieval success.⁴ Despite the large number of studies that have documented and characterized SMEs across a wide range of memory tasks and encoding manipulations, the relative contributions of endogenous and external factors have yet to be established.

Free recall studies of SMEs typically compare brain activity associated with the encoding of subsequently recalled and non-recalled items within a given list. Some of the strongest predictors of recall performance are characteristics of individual items (e.g., their pre-experimental familiarity or their position in the study list).⁵⁻⁷ Such idiosyncratic item-level effects are therefore serious confounds in item-level SME analyses and difficult to control, because repetition of items across lists would produce carry-over effects.

To limit these item-level effects in our examination of broader external factors that also affect recall performance (such as session-level time-of-day effects or list-level proactive interference effects), we computed list-level SMEs by averaging the epochs of brain activity following the presentation of individual study items across study lists. Specifically, we analyzed EEG recordings

from 97 individuals who each studied and recalled 24 word lists in each of at least 20 experimental sessions that took place over the course of several weeks. We trained ridge regression models to predict the (logit-transformed) proportion of recalled items for each list, $p(\text{rec})$, on the basis of spectral EEG features that we averaged over recordings during all encoding periods in that list. Additionally, we leveraged a prior statistical model of memory performance which identified several critical variables predicting recall performance across both lists and sessions.⁴ By removing linear effects of these variables, we uncovered the components of neural activity that predict the residual recallability of studied items. Comparing SMEs for these residuals with those obtained for raw recall performance thus allowed us to estimate the relative contributions of endogenous neural variability and external factors to the SME. Throughout this paper we assessed our ability to predict recall performance with a leave-one-session-out cross-validation procedure (see methods for details).

Results

Figure 1 shows the mean proportion of recall as the function of several external variables that affect recall performance for entire sessions (intersession predictors) and for individual lists within each session (interlist predictors). Specifically, we considered sleep duration in the night prior to the free recall test, time of day, and self-rated alertness at the beginning of the experimental session as intersession predictors and experimental block within each session, the list number within each block, and the average “recallability” of items within each list as interlist predictors.⁴ We are showing the effects of these variables across all participants (discretized into two bins for each of the

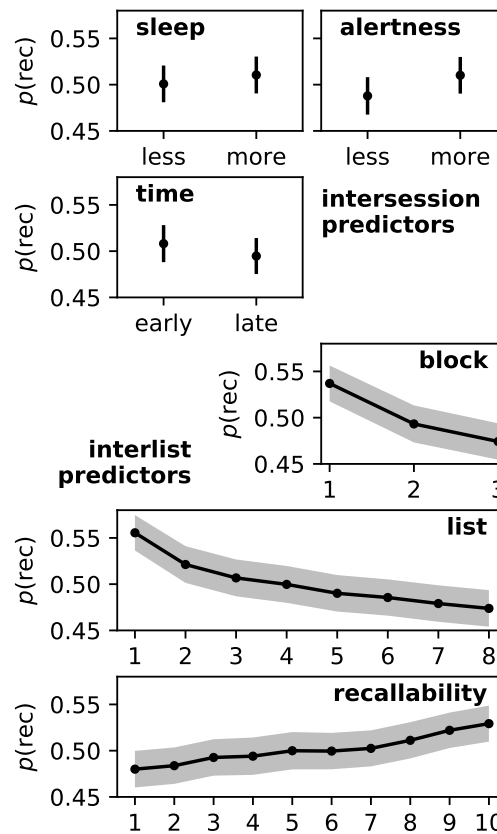


Figure 1: Mean probability of recall (and associated 95% confidence intervals) as a function of intersession (amount of sleep, rated alertness, and time of day) and interlist (block number within a session, list number within a block, and mean recallability of items within a list) predictors. For the purpose of this visualization we discretized each individual's intersession predictors into two bins and mean recallability scores into ten bins, but our analyses applied separately to the full data set from each individual.

intersession predictors and into ten bins for recallability) for illustrative purposes, but we applied all of our analyses separately to the full data from each individual. Additionally, we also considered the effect of session number (which was heterogeneous across participants with some showing increased performance with increasing practice and some showing a decline in performance) as an additional predictor in our intersession and interlist regression models (described below). Detailed analyses of the effects of these variables on recall performance in a large subset of this data set were the focus of a previous study.⁴

Figure 2A illustrates brain activity associated with the study of individual items from two adjacent lists that were associated with relatively high (21/24) and relatively low (11/24) proportions of recall respectively (to accommodate size constraints we omitted the middle 12 items from each list in this figure). The brain activity shown in these time-frequency plots extends beyond the time range for the epochs we used in our analyses (0.3–1.6 s after study-word onset) to 0.8 s before each item’s onset to capture most of the brain activity during the presented sections of the study lists (additional variable-duration inter-stimulus intervals are not shown). The predicted probability that an item will subsequently be recalled according to an item-based classifier (see methods) is shown below each sub-panel.

Figure 2B illustrates brain activity across encoding periods for subsequently recalled and subsequently unrecalled encoding epochs across all participants in our study. This panel shows distinct patterns of brain activity as a function of subsequent recall — a standard subsequent memory effect. Given the strong effects of item-level characteristics on recall performance, it is possible

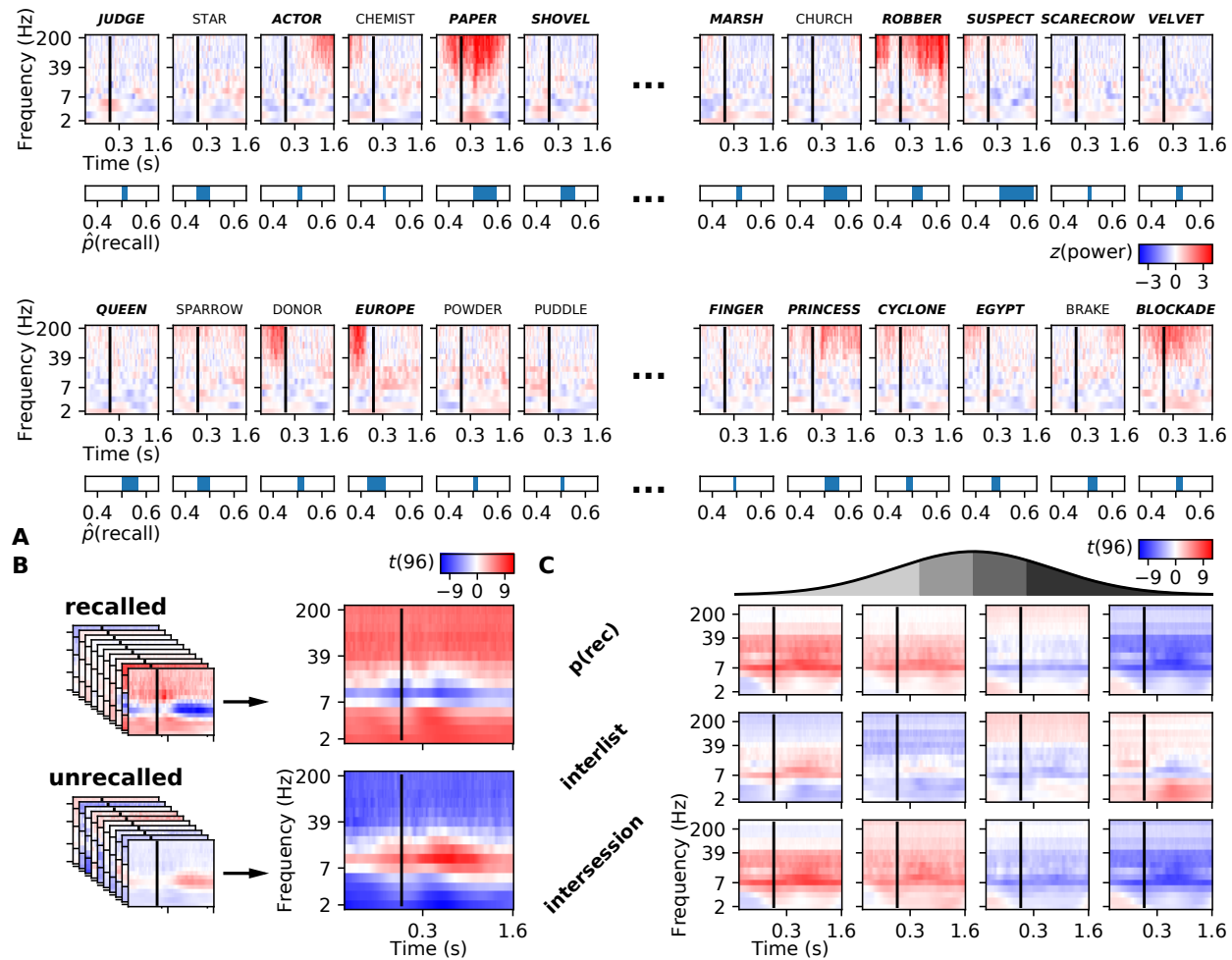


Figure 2: z -transformed power around the presentation of study words during one participant's (ID: 374) 16th experimental session. The two rows show excerpts from the 4th and 5th study lists in that session. The study words are indicated at the top of each sub-panel with bold italic font indicating subsequent recall. The horizontal bar graphs at the bottom of each sub-panel show the output of the logistic regression classifier (see methods). **B:** Average power for subsequently recalled (top) and subsequently unrecalled (bottom) words during study across all participants. **C:** Average power for quartiles of $p(\text{rec})$ (top) as well as interlist (middle) and intersession (bottom) residuals across all participants. For this visualization, we aggregated activity across electrodes in four superior regions of interest (illustrated in Figure 4). All our analyses were based on data from individual electrodes with no discretization into quantiles. Vertical black lines indicate word onset. The time range for our encoding epochs (0.3–1.6 s after study-word onset) are labeled on the time axes.

that they explain a large proportion of the variance in item-level SMEs. Additionally, it is possible that any endogenous variability driving SMEs is relatively fast, varying on the order of seconds (i.e., the time devoted to the study of individual items in typical memory experiments) rather than tens of seconds (i.e., the time encompassing a full study list) or longer. It was therefore not clear that brain activity averaged over the individual study periods would be similarly informative about list-level recall performance as standard item-level SMEs. For our list-level analyses, we averaged encoding epochs of brain activity within each list to predict (logit transformed) $p(\text{rec})$ and residuals from the interlist and intersession regression models. To illustrate how brain activity covaried with these continuous variables, we partitioned lists into quartiles based on the (residuals of) recall performance and show the average brain activity for lists in each quartile across all participants in Figure 2C (all of our regression analyses described below are based on the continuous measures and recordings from individual electrodes). Figure 2C suggests that list-level brain activity varies considerably as a function of list-level recall performance and that this effect is largely preserved even when accounting for the linear effects of external factors, especially those varying between sessions (intersession predictors). Furthermore, a comparison of Panels B and C reveals apparent similarities between neural features predicting item-level and list-level recall performance with a decrease in α power (around 10 Hz) associated with better recall performance. We confirmed these impressions with our detailed multivariate analyses to which we now turn.

To compare the sizes of our list-level SME to the classic item-level SME, we trained an L2 penalized logistic regression (LR) model to predict subsequent recall of individual items (again using a leave-one-session-out cross-validation procedure to measure classification performance;

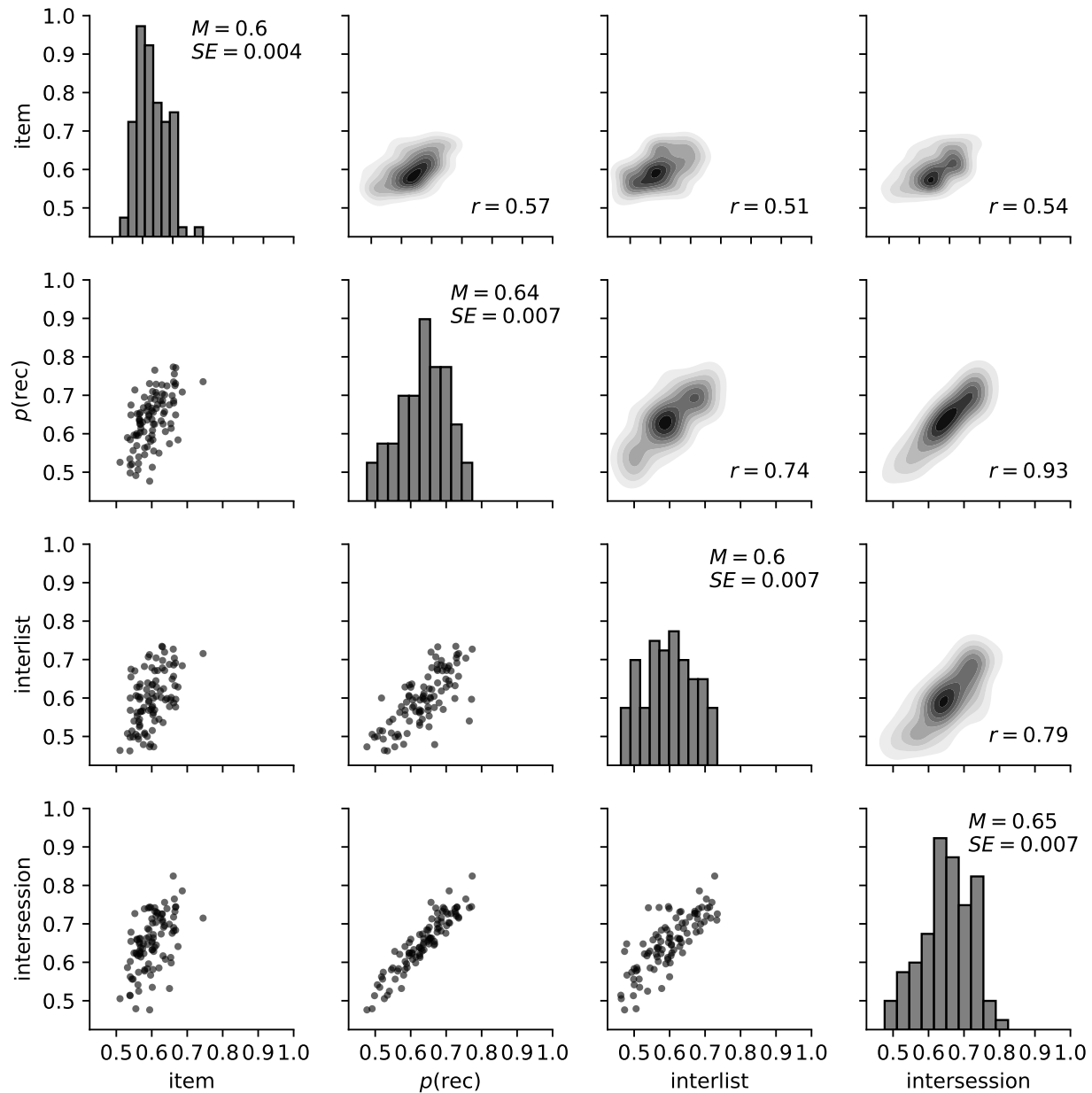


Figure 3: Areas under the ROC functions (AUCs) for classifier performance predicting subsequent memory for individual items (using a logistic regression classifier; item), lists of items (predicting the probability of recall across list items; $p(\text{rec})$), as well as for residuals of list-level recall performance after regressing out interlist and intersession predictors. The lower triangle shows scatter plots for each pair of AUCs across participants. The upper triangle shows bivariate kernel density estimates of these same data with the corresponding correlations. The main diagonal shows histograms of each classifier's AUCs with the corresponding means and standard errors.

example outputs are shown in Figure 2A). For classification problems, the area under the receiver operating characteristic (ROC) function (AUC) provides a convenient index of classification performance with an AUC of 0.5 corresponding to chance performance and an AUC of 1.0 indexing perfect classification.⁸ To allow direct comparisons between the performance of the item-level classifier and our ridge regression models predicting $p(\text{rec})$ we also calculated AUCs for our regression models by discretizing the proportion of list-level recalls. Specifically, for these analyses, we treated lists whose $p(\text{rec})$ exceeded the total proportion of recalled items in a session as the target category and all other lists as the non-target category. Figure 3 shows AUCs for the item-level classifier as well as for three different list-level regression models which we will discuss in turn.

The list-level regression model predicting $p(\text{rec})$ yielded a mean AUC of 0.64 which was significantly higher than that for the item-level LR classifier ($M = 0.6$; $t(96) = 6.168$, $SE = 0.006$, $p < 0.001$). This demonstrates that spectral features averaged over encoding periods effectively predict list-level recall performance. This pattern of results is consistent with the idea that slow (on the order of tens of seconds) endogenous processes rather than idiosyncratic item-level characteristics or fast (on the order of seconds) endogenous variation are the primary drivers of the SME. It is difficult, however, to directly compare list-level and item-level results because the averaging over 24 individual encoding epochs for the list-level analyses or the corresponding 24-fold larger training set for the item-level analysis could each affect the relative performance of the respective analysis. To assess the effect of averaging, we also applied the item-level LR classifier to the average brain activity of all recalled and all unrecalled trials within each list (i.e., at most 2 in-

stead of 24 observations per list; not shown in the figure). This averaging yielded a mean AUC of 0.68 which was significantly larger than that for the standard item-level classifier ($t(96) = 15.294$, $SE = 0.005$) and the regression model predicting $p(\text{rec})$ ($t(96) = 15.294$, $SE = 0.005$, both $ps < 0.001$). The fact that the advantage of averaging outweighed any detriment of reducing the number of training samples is also consistent with relevant signals varying on slow (on the order of tens of seconds) time scales, especially given that there are strong sequential dependencies in recall performance⁹ (resulting in the preferential averaging of adjacent items in the averaged item-level LR classifier).

Given that item-level characteristics and fast endogenous variation did not appear to substantially contribute to the SME, we next considered the extent to which external variables affecting recall performance for entire sessions (intersession predictors: sleep, alertness, and time of day) and those that affect recall performance at the list-level (interlist predictors: block, list, recallability)⁴ are driving differences in brain activity that predict recall success. To the extent that either set of variables can explain the SME, we can conclude that it also does not reflect slow endogenous variability on the order of tens of seconds (i.e., lists) or days (i.e., sessions). We constructed interlist and intersession regression models (both models also included session number as a predictor) to remove linear effects of the respective external variables on $p(\text{rec})$. We then predicted the resulting list-level residuals with ridge regression models using the same spectral EEG features as for our list-level regression model predicting $p(\text{rec})$. Figure 3 shows that AUCs for the interlist and intersession regression models respectively matched ($M = 0.6$) or exceeded ($M = 0.65$, $t(96) = 8.354$, $SE = 0.006$, $p < 0.001$) those for the item level classifier, demonstrating that spectral features

effectively predict list-level performance even after accounting for linear effects from several external variables that affect recall. These results thus rule out these factors as major contributors to the SME, suggesting that SMEs predominantly reflect slow endogenous variability in cognitive function.

Whereas the AUCs for the interlist regression models were significantly lower than those of the regression models predicting $p(\text{rec})$ ($t(96) = 7.389$, $SE = 0.005$, $p < 0.001$), the AUCs for the intersession regression models exceeded those for the other list-level regression models ($t(96) = 5.828$ and 11.464 , $SE = 0.003$ and 0.005 , $ps < 0.001$, for comparisons with the $p(\text{rec})$ and interlist models, respectively; Figure 3). This pattern of results indicates some effects of interlist factors on our measures of brain activity predicting recall performance, leading to a reduction in model performance when linear effects of interlist predictors were removed. The fact that the intersession models were better able to generalize across sessions indicates that relevant brain activity varying across sessions was not effectively captured by our models (because we used a leave-one-session-out cross-validation procedure to measure model performance, AUCs index the ability of our models to generalize across sessions). Thus, removing linear effects of intersession predictors removed variability that the models could not account for, leading to increased performance. These results establish a small role for list-level effects due to external factors (e.g., proactive interference) in the SME in addition to strong effects of endogenous variability in encoding processes.

Figure 3 also highlights substantial correlations between AUCs for the different models.

This suggests that the different models use brain activity similarly to predict (residuals of) recall performance. It is difficult, however, to interpret the levels of these correlations in light of the fact that the dependent measures also correlate substantially—a previous analysis⁴ showed a reduction of variability of the residuals for the interlist and intersession models relative to $p(\text{rec})$ of only around 11% on average, leaving most of the variability in recall performance unaccounted for by external variables.

A standard measure of performance for regression models is the correlation between predicted and actual values of the dependent measures. These correlations mirror the pattern of the AUCs shown in Figure 3 with $r = 0.26, 0.28$, and 0.19 for $p(\text{rec})$, intersession residuals, and interlist residuals respectively (all pairwise differences were statistically significant, $t(96) = 8.370\text{--}11.695$, $SE = 0.003\text{--}0.008$, $ps < 0.001$). The point-biserial correlation between predictions from the item-level classifier and recall status of individual items was 0.16 (0.30 for the averaging item-level classifier). This confirms the above AUC-based analyses indicating the effectiveness of spectral features in predicting list-level performance and the ability of our models to capture some brain activity associated with interlist, but not intersession, predictors (because of the better performance for the intersession models and the reduced performance of the interlist models relative to the models predicting $p(\text{rec})$ as explained above).

In addition to investigating the correlations between predictions from the different regression models and the corresponding dependent measures, we can also assess the extent to which the different models generalize to predicting the other measures.¹ This analysis reveals an advantage

¹This is conceptually similar to a cross-decoding approach where models trained on one data set are used for

for models trained on intersession residuals, even when these were tested on $p(\text{rec})$ or interlist residuals. To assess the size of these differences, we removed the linear effects of the measure each model was trained on from the generalization measures and computed the (semi-partial) correlations between the model predictions and the resulting residuals. The semi-partial correlations between predictions of models trained on intersession-residuals and the other two measures were positive ($M = 0.1$ for both $p(\text{rec})$ and interlist-residuals; $t(96) = 17.43$ and 13.339 , $SE = 0.006$ and 0.008 , $ps < 0.001$, respectively). This confirms that the performance advantage for models trained on intersession residuals generalizes to the prediction of $p(\text{rec})$ and interlist residuals—a result that complements the above finding suggesting that removing linear effects of intersession predictors eliminates variability in recall performance that is not effectively captured by our measures of brain activity. In contrast, the semi-partial correlations between predictions of models trained on interlist-residuals and the other two measures were negative ($M = -0.18$ and -0.23 , $t(96) = 9.463$ and 14.671 , $SE = 0.019$ and 0.016 , for the $p(\text{rec})$ and intersession residuals respectively, both $ps < 0.001$). This indicates that the relative disadvantage for models trained on interlist residuals generalizes to the prediction of $p(\text{rec})$ and intersession residuals, consistent with our measures of brain activity being sensitive to interlist variables. The only other semi-partial correlations significantly deviating from 0 were those between predictions of the models trained on $p(\text{rec})$ and the interlist-residuals ($M = 0.07$, $t(96) = 9.831$, $SE = 0.007$, $p < 0.001$), reflecting the fact that models trained on $p(\text{rec})$ were better able to capitalize on brain activity that is relevant for predicting recall performance than models that could not make use of brain activity that reflects

predictions on a different data set.⁹ In the current application we train models on identical features to predict different measures of recall performance rather than predicting the same dependent measure in different data sets.

interlist predictors.

Figure 3 showed high correlations between the performances of the different models predicting item and list-level recall which suggests that there is considerable overlap between the patterns of brain activity predicting these measures. We investigated this relationship by correlating power across a range of frequencies and regions of interest (ROIs) with each of the measures of recall performance. These correlations exhibited a consistent pattern with low (negative) correlations in the θ and α range ($\approx 5\text{--}10$ Hz) which increased for higher (and lower) frequencies (Figure 4). For the (point-biserial) correlation of brain activity with item-level recall, we observed negative correlations in the θ and α range and positive correlations in the γ (> 30 Hz) range, consistent with numerous findings showing that decreased power in lower frequencies and increased power in higher frequencies predicts subsequent memory.^{9–12} As shown in Figure 4, the correlations for the list-level measures of recall performance exhibited qualitatively very similar patterns, confirming that the different ways of calculating SMEs leverage brain activity in similar ways (see also Figure 2B and C).

The similarity in how brain activity correlates with different measures of recall performance complements our analysis of correlations between AUCs associated with different regression models (Figure 3). Just like that analysis, however, this similarity is difficult to interpret in light of substantial correlations between the dependent measures. To directly assess how brain activity covaries with variability that is specific to intersession and interlist predictors (removing linear effects of $p(\text{rec})$), we therefore correlated brain activity with corresponding residuals (intersession| $p(\text{rec})$)

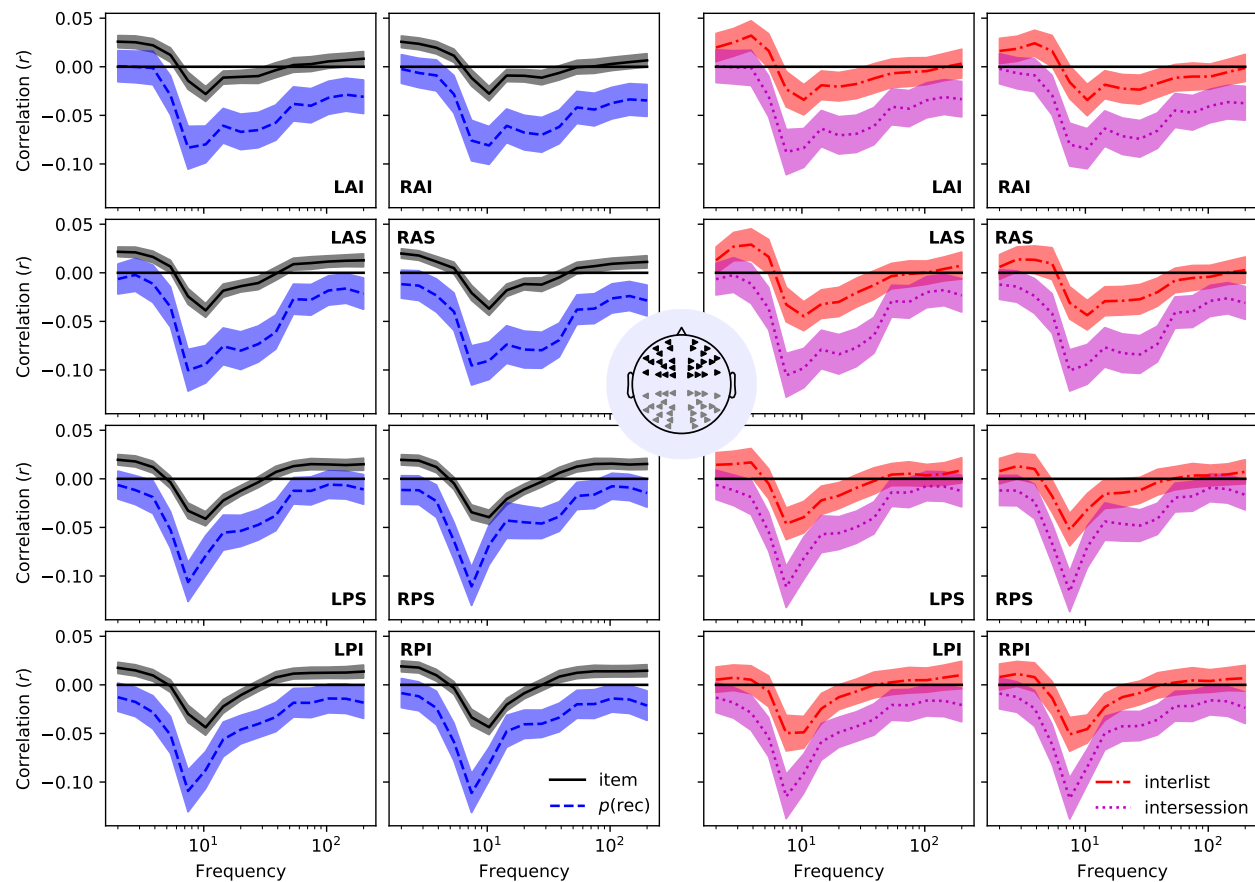


Figure 4: Correlations between mean power in each frequency across electrodes within each region of interest (ROI) and measures of recall performance (recall of individual items, $p(\text{rec})$, and residuals from the interlist and intersession models). The inset in the middle of the figure illustrates the locations of the ROIs and each panel includes an ROI label with the first letter indicating the hemisphere (L: left, R: right), the second letter distinguishing between anterior (A) and posterior (P) ROIs, and the last letter specifying the ROI position as either inferior (I) or superior (S). Zero is indicated as are 95% confidence intervals (shaded regions).

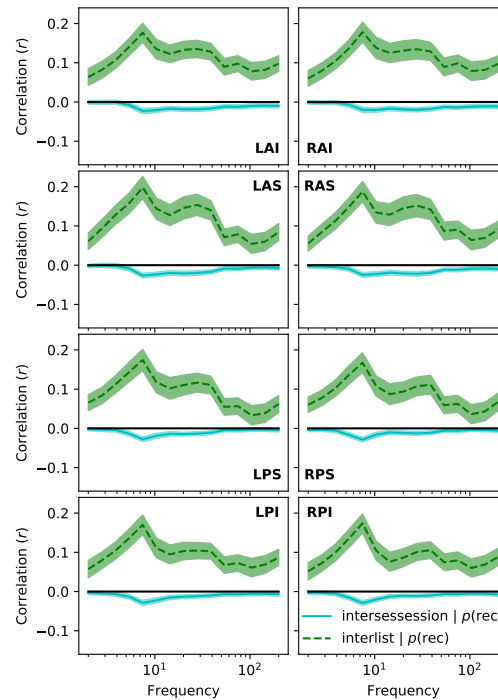


Figure 5: Correlations between mean power in each frequency across electrodes within each region of interest (ROI) and intersession and interlist residuals after regressing out linear effects of $p(\text{rec})$ (intersession| $p(\text{rec})$ and interlist| $p(\text{rec})$ respectively). Each panel shows these correlations for a different ROI (labeled and arranged as in Figure 4). Zero is indicated as are 95% confidence intervals (shaded regions).

and interlist| $p(\text{rec})$, respectively; Figure 5). As is evident from Figure 5, correlations of brain activity with intersession| $p(\text{rec})$ residuals were close to zero and varied little across frequencies or ROIs, consistent with the above analyses indicating that our measures of brain activity did not capture much of the variability in recall performance associated with intersession predictors. The correlations of brain activity with interlist| $p(\text{rec})$, however, were relatively strong, complementing the above analyses indicating that our measures of brain activity are sensitive to interlist predictors of recall performance.

Discussion

Whether and how a studied item is encoded and subsequently retrieved during a free recall task is, by design, not subject to complete experimental control. Indeed, recalled and not-recalled items tend to differ on a number of dimensions. Prior work has shown that neural activity just before the presentation of individual items predicts subsequent memory performance, demonstrating SMEs that are independent of specific item characteristics.^{13–16} Nevertheless, task-related variables also strongly predict memory performance and could be driving SMEs even when they are not linked to specific item characteristics (e.g., recalled items tend to disproportionately come from early list positions, a “primacy” effect).⁷ Thus, any comparison of brain activity during the study of items as a function of their subsequent recall is fraught with confounds, complicating the interpretation of the diagnostic neural signals. We avoided some of these confounds by assessing list-level SMEs, aggregating brain activity across the study periods of all items within a list to predict list-level recall. Our demonstration of substantial list-level SMEs (Figure 3) and similar predictive patterns

of brain activity for item and list-level SMEs (Figures 2 and 4), shows that item-level SMEs are not mainly driven by external variables differentiating items within a study list. This result also suggests the presence of endogenous neural variation at slow time scales (items in a study list were presented over the course of about a minute) that predicts subsequent memory.

Even when aggregating across items within a list, a range of confounding variables remain. By studying 97 individuals who each participated in up to 23 experimental sessions, we were able to model the effects of several external variables that affect list-level recall performance. This enabled us to not only relate brain activity to the proportion of recalled items in each list, but also to residuals of recall performance after accounting for effects of these external variables. Following previous work,⁴ we partitioned these external variables into those that varied across lists (interlist) and those that varied across sessions (intersession). Accounting for interlist variables reduced the list-level SME slightly (Figures 3 and 4). This suggests that some, but not all, of the list-level SME reflects the effects of interlist variables. Accounting for intersession variables, on the other hand, slightly increased the size of the SME, demonstrating that the list-level SME does not include substantial contributions from these variables (Figures 3 and 4; see also Figure 5).

Distinguishing between effects of external variables and endogenous processes is notoriously difficult, because it is impossible to control for effects of all possible external factors. Additionally some external factors (e.g., drug consumption or exercise) can have long-lasting and/or variable effects, making it difficult to establish their relationship with behavior. Indeed, the distinction between external and endogenous effects can be blurry, especially when external variables (such

as time of day) correlate with endogenous processes (e.g., physiological changes due to circadian rhythms). In our investigation of variability in recall performance, we controlled for the major variables known to affect episodic memory. We also considered broad variables (such as recallability, time of day, and alertness) that were meant to capture the joint effects of large sets of more specific variables (e.g., features of the individual words within a study list, number of waking hours, or effects of caffeine consumption). Thus, we believe that the joint effects of external variables beyond those considered as predictors in our interlist and intersession models are likely to be too small to account for a substantial fraction of the remaining variability in recall performance or the SME.

When we controlled for the effects of sleep, alertness, and time of day, our ability to predict list-level recall from brain activity increased. This indicates that these variables did not substantially contribute to the list-level SME we observed (and hence removing their effects improved generalization of our models). Our results thus highlight the need to distinguish between variables that affect recall performance and those whose effects manifest in our measures of brain activity. Considering additional variables that affect recall performance therefore need not reduce our estimate of the contributions of endogenous factors to the SME.

The fact that substantial SMEs remained after accounting for a comprehensive set of external variables may appear in conflict with findings that task context can affect the specific form of SMEs, at least for recognition memory.^{15,17–20} Task context manipulations in these studies were designed to directly affect encoding processes (e.g., by asking participants to perform different tasks on the study items) and their effects on SMEs suggest that they were successful. Here we

show that in the absence of direct manipulations of how study items are presented or processed, external variables do not substantially contribute to the SME even when they predict subsequent recall. These findings indicate that SMEs are not only effective measures of memory formation, but that they reflect endogenous states that drive variability in cognitive function.

Our findings align well with reports of sequential dependencies in human performance^{4,21–23} as well as with those of slow endogenous neural fluctuations that drive variability in evoked brain activity and overt behavior.^{24–30} Previous investigations of endogenous variability in neural activity and performance have relied on exact repetitions of stimuli across many experimental trials to limit variability in external factors. In order to study the effects of endogenous variability on recall performance, we took a complementary approach by statistically removing the effects of a comprehensive set of external factors. Despite the differences in methodologies and tasks, the conclusions are remarkably consistent in establishing an important role for slowly varying fluctuations in neural activity as drivers of variability in human cognition.

Methods

Participants We analyzed data from 97 young adults (18–35) who completed at least 20 sessions in Experiment 4 of the Penn Electrophysiology of Encoding and Retrieval Study (PEERS) in exchange for monetary compensation. This study was approved by the Institutional Review Board at the University of Pennsylvania and we obtained informed consent from all participants. Recall performance for a large subset of the current data set was previously reported,⁴ but this is the first report of electrophysiological data from this experiment. Data from PEERS experiments are freely

available at <http://memory.psych.upenn.edu> and have been reported in several previous publications.^{31–38} Our analyses included data from all participants with at least 20 sessions.

Experimental task Each of up to 23 experimental sessions consisted of 24 study lists that each were followed by a delayed free recall test. Specifically, each study list presented 24 session-unique English words sequentially for 1,600 ms each with a blank inter-stimulus interval that was randomly jittered (following a uniform distribution) between 800 and 1,200 ms. After the last word in each list, participants were asked to solve a series of arithmetic problems of the form $A + B + C = ?$ where, A , B , and C were integers in $[1, 9]$. Participants responded to each problem by typing the result and were rewarded with a monetary bonus for each correctly solved equation. These arithmetic problems were displayed until 24 s had elapsed and were then followed by a blank screen randomly jittered (following a uniform distribution) to last between 1,200 and 1,400 ms. Following this delay, a row of asterisks and a tone signaled the beginning of a 75 s free recall period. A random half of the study lists (except for the first list in each session) were also preceded by the same arithmetic distractor task which was separated from the first study-item presentation by a random delay jittered (following a uniform distribution) to last between 800 and 1,200 ms. Each session was partitioned into 3 blocks of 8 lists each and blocks were separated by short (approximately 5 min) breaks. At each session participants were asked to rate their alertness and indicate the number of hours they had slept in the previous night.

Stimuli Across all lists in each session the same 576 common English words (24 words in each of 24 lists) were presented for study, but their arrangement into lists differed from session to session (subject to constraints on semantic similarity³¹). These 576 words were selected from a larger

word pool (comprising 1,638 words) used in other PEERS experiments. The 576-word subset of this pool used in the current experiment were selected to maximize homogeneity, by removing words that were atypical in frequency, concreteness, or emotional valence. Many participants also returned for a 24th session that used words from the entire 1,638-word pool, but we are not reporting data from that session here. We estimated the mean recallability of items in a list from the proportion of times each word within the list was recalled by other participants in this study.

EEG data collection and processing Electroencephalogram (EEG) data were recorded with either a 129 channel Geodesic Sensor net using the Netstation acquisition environment (Electrical Geodesics, Inc.; EGI) or with a 128 channel Biosemi Active Two system. EEG recordings were re-referenced offline to the average reference. Because our regression models weighted neural features with respect to their ability to predict (residuals of) recall performance in held out sessions, we did not try to separately eliminate artifacts in our EEG data. Data from each participant were recorded with the same EEG system throughout all sessions and for those sessions recorded with the Geodesic Sensor net, we excluded 26 electrodes that were placed on the face and neck, rather than the scalp, from further analyses. The EGI system recorded data with a 0.1 Hz high-pass filter and we applied a corresponding high-pass filter to the data collected with the Biosemi system. We used MNE,^{39,40} the Python Time-Series Analysis (PTSA) library (https://github.com/pennmem/ptsa_new), Sklearn⁴¹ and custom code for all analyses.

We first partitioned EEG data into epochs starting 800 ms before the onset of each word in the study lists and ending with its offset (i.e., 1,600 ms after word onset). We also included an additional 1,200 ms buffer on each end of each epoch to eliminate edge effects in the wavelet

transform. We calculated power in 15 logarithmically spaced frequencies between 2 and 200 Hz, applied a log-transform, and down-sampled the resulting time series of log-power values to 50 Hz. We then truncated each epoch to 300–1,600 ms after word onset. For the item-based classifier we used each item’s mean power in each frequency across this 1,300 ms interval as features to predict subsequent recall (we also present result for an averaged item-based classifier that aggregated these intervals across all subsequently recalled and all subsequently unrecalled items in each list). For the list-based regression models we averaged these values across all items in each list to predict (residuals of) list-level recall.

For the analyses shown in Figures 4 and 5, we partitioned electrodes into the 6 regions of interest (ROIs) illustrated in Figure 4 (we also aggregated over electrodes in the 4 superior ROIs in Figure 2). This choice of ROIs follows a range of studies that used these or very similar ROIs to characterize the spatial distribution of EEG effects.⁴² All of our classification and regression models, however, used measures from individual electrodes as input without any averaging into ROIs.

Item-based classifier For the item-based classifier we used a nested cross-validation procedure to simultaneously determine the regularization parameter and performance of L2-regularized logistic regression models predicting each item’s subsequent recall. At the top level of the nested cross-validation procedure we held out each session once—these held out sessions were used to assess the performance of the models. Within the remaining sessions, we again held out each session once—these held-out sessions from within each top-level cross-validation fold were used to determine the optimal regularization parameter, C , for Sklearn’s LogisticRegression class. We

fit models with 9 different C values between 0.00002 and 1 to the remaining sessions within each cross-validation fold and evaluated their performance as a function of C on the basis of the held out sessions within this fold. We then fit another logistic regression model using the best-performing C value to all sessions within each cross-validation fold and determined the model predictions on the sessions that were held-out at the top level. We calculated the area under the ROC function on the basis of the predictions from these held-out sessions.

List-based regression models For the list-based regression models we followed the same procedure as for the item-based classifier to determine the optimal level of regularization for ridge regression models predicting (residuals of) list-level recall performance. Specifically, we used the same nested cross-validation procedure described above to determine optimal values for α (corresponding to $1/C$), the regularization parameter in Sklearn's Ridge class, testing 9 values between 1 and 65536. We applied these models to the (logit-transformed) proportion of items recalled for each list, $p(\text{rec})$, as well as to the residuals from the interlist and intersession models as described in the results section.⁴

Data availability Data from this experiment are freely available at <http://memory.psych.upenn.edu>.

Code availability Data analysis code from this manuscript is freely available at <http://memory.psych.upenn.edu>.

1. Paller, K. A. & Wagner, A. D. Observing the transformation of experience into memory.
Trends in Cognitive Sciences **6**, 93–102 (2002).
2. Kim, H. Neural activity that predicts subsequent memory and forgetting: A meta-analysis of
74 fMRI studies. *NeuroImage* **54**, 2446–2461 (2011).
3. Hanslmayr, S. & Staudigl, T. How brain oscillations form memories — a processing based
perspective on oscillatory subsequent memory effects. *NeuroImage* **85**, 648–655 (2014).
4. Kahana, M. J., Aggarwal, E. V. & Phan, T. D. The variability puzzle in human memory.
Journal of Experimental Psychology: Learning, Memory, and Cognition (2018).
5. DeLosh, E. L. & McDaniel, M. A. The role of order information in free recall: Application
to the word-frequency effect. *Journal of Experimental Psychology: Learning, Memory, and
Cognition* **22**, 1136–1146 (1996).
6. Merritt, P. S., DeLosh, E. L. & McDaniel, M. A. Effects of word frequency on individual-
item and serial order retention: Tests of the order-encoding view. *Memory & Cognition* **34**,
1615–1627 (2006).
7. Murdock, B. B., Jr. The serial position effect of free recall. *Journal of Experimental Psychol-
ogy* **64**, 482–488 (1962).
8. Fawcett, T. An introduction to ROC analysis. *Pattern Recognition Letters* **27**, 861–874 (2006).

- 403 9. Weidemann, C. T. *et al.* Neural activity reveals interactions between episodic and semantic
404 memory systems during retrieval. *Journal of Experimental Psychology: General* **148**, 1–12
405 (2019).
- 406 10. Hanslmayr, S., Staudigl, T. & Fellner, M.-C. Oscillatory power decreases and long-term mem-
407 ory: the information via desynchronization hypothesis. *Frontiers in Human Neuroscience* **6**
408 (2012).
- 409 11. Burke, J. F. *et al.* Human intracranial high-frequency activity maps episodic memory forma-
410 tion in space and time. *NeuroImage* **85 Pt. 2**, 834–843 (2014).
- 411 12. Long, N. M. & Kahana, M. J. Successful memory formation is driven by contextual encoding
412 in the core memory network. *NeuroImage* **119**, 332–337 (2015).
- 413 13. Sweeney-Reed, C. M. *et al.* Pre-stimulus thalamic theta power predicts human memory for-
414 mation. *NeuroImage* **138**, 100–108 (2016).
- 415 14. Otten, L. J., Quayle, A. H., Akram, S., Ditewig, T. A. & Rugg, M. D. Brain activity before an
416 event predicts later recollection. *Nature Neuroscience* **9**, 489–491 (2006).
- 417 15. Fellner, M.-C., Bäuml, K.-H. T. & Hanslmayr, S. Brain oscillatory subsequent memory effects
418 differ in power and long-range synchronization between semantic and survival processing.
419 *NeuroImage* **79**, 361–370 (2013).
- 420 16. Guderian, S., Schott, B. H., Richardson-Klavehn, A. & Düzel, E. Medial temporal theta state
421 before an event predicts episodic encoding success in humans. *Proceedings of the National*
422 *Academy of Sciences* **106**, 5365–5370 (2009).

- 423 17. Kamp, S.-M., Bader, R. & Mecklinger, A. ERP subsequent memory effects differ between
424 inter-item and unitization encoding tasks. *Frontiers in Human Neuroscience* **11** (2017).
- 425 18. Summerfield, C. & Mangels, J. A. Dissociable neural mechanisms for encoding predictable
426 and unpredictable events. *Journal of Cognitive Neuroscience* **18**, 1120–1132 (2006).
- 427 19. Otten, L. J. & Rugg, M. D. Electrophysiological correlates of memory encoding are task-
428 dependent. *Cognitive Brain Research* **12**, 11–18 (2001).
- 429 20. Staudigl, T. & Hanslmayr, S. Theta oscillations at encoding mediate the context-dependent
430 nature of human episodic memory. *Current Biology* **23**, 1101–1106 (2013).
- 431 21. Gilden, D., Thornton, T. & Mallon, M. 1/f noise in human cognition. *Science* **267**, 1837–1839
432 (1995).
- 433 22. Mueller, S. T. & Weidemann, C. T. Decision noise: An explanation for observed violations of
434 signal detection theory. *Psychonomic Bulletin & Review* **15**, 465–494 (2008).
- 435 23. Verplanck, W. S., Collier, G. H. & Cotton, J. W. Nonindependence of successive responses
436 in measurements of the visual threshold. *Journal of Experimental Psychology* **44**, 273–282
437 (1952).
- 438 24. Monto, S., Palva, S., Voipio, J. & Palva, J. M. Very slow EEG fluctuations predict the dynamics
439 of stimulus detection and oscillation amplitudes in humans. *Journal of Neuroscience* **28**, 8268–
440 8272 (2008).

- 441 25. Schroeder, C. E. & Lakatos, P. Low-frequency neuronal oscillations as instruments of sensory
442 selection. *Trends in Neurosciences* **32**, 9–18 (2009).
- 443 26. Arieli, A., Sterkin, A., Grinvald, A. & Aertsen, A. Dynamics of ongoing activity: Explanation
444 of the large variability in evoked cortical responses. *Science* **273**, 1868–1871 (1996).
- 445 27. Fox, M. D., Snyder, A. Z., Zacks, J. M. & Raichle, M. E. Coherent spontaneous activity
446 accounts for trial-to-trial variability in human evoked brain responses. *Nature Neuroscience* **9**,
447 23–25 (2005).
- 448 28. Fox, M. D., Snyder, A. Z., Vincent, J. L. & Raichle, M. E. Intrinsic fluctuations within cortical
449 systems account for intertrial variability in human behavior. *Neuron* **56**, 171–184 (2007).
- 450 29. Fox, M. D. & Raichle, M. E. Spontaneous fluctuations in brain activity observed with func-
451 tional magnetic resonance imaging. *Nature Reviews Neuroscience* **8**, 700–711 (2007).
- 452 30. Raichle, M. E. The restless brain: how intrinsic activity organizes brain function. *Philosophi-
453 cal Transactions of the Royal Society B: Biological Sciences* **370**, 20140172 (2015).
- 454 31. Healey, M. K., Crutchley, P. & Kahana, M. J. Individual differences in memory search and
455 their relation to intelligence. *Journal of Experimental Psychology: General* **143**, 1553–1569
456 (2014).
- 457 32. Healey, M. K. & Kahana, M. J. Is memory search governed by universal principles or idiosyn-
458 cratic strategies? *Journal of Experimental Psychology: General* **143**, 575–596 (2014).

33. Healey, M. K. & Kahana, M. J. Age-related changes in the dynamics of memory encoding processes provide a biomarker of successful aging. *Manuscript Submitted for publication* (2018).
34. Lohnas, L. J. & Kahana, M. J. Parametric effects of word frequency in memory for mixed frequency lists. *Journal of Experimental Psychology: Learning, Memory, and Cognition* **39**, 1943–1946 (2013).
35. Siegel, L. L. & Kahana, M. J. A retrieved context account of spacing and repetition effects in free recall. *Journal of Experimental Psychology: Learning, Memory, and Cognition* **40**, 755–764 (2014).
36. Lohnas, L. J., Polyn, S. M. & Kahana, M. J. Expanding the scope of memory search: Modeling intralist and interlist effects in free recall. *Psychological Review* **122**, 337–363 (2015).
37. Weidemann, C. T. & Kahana, M. J. Assessing recognition memory using confidence ratings and response times. *Royal Society Open Science* **3**, 150670 (2016).
38. Weidemann, C. T. & Kahana, M. J. Dynamics of brain activity reveal a unitary recognition signal. *Journal of Experimental Psychology: Learning, Memory, and Cognition* **45**, 440–451 (2019).
39. Gramfort, A. *et al.* MEG and EEG data analysis with MNE-python. *Frontiers in Neuroscience* **7**, 267 (2013).
40. Gramfort, A. *et al.* MNE software for processing MEG and EEG data. *NeuroImage* **86**, 446–460 (2014).

479 41. Pedregosa, F. *et al.* Scikit-learn: Machine learning in Python. *Journal of Machine Learning*
480 *Research* **12**, 2825–2830 (2011).

481 42. Weidemann, C. T., Mollison, M. V. & Kahana, M. J. Electrophysiological correlates of
482 high-level perception during spatial navigation. *Psychonomic Bulletin & Review* **16**, 313–319
483 (2009).

484 **Acknowledgements** his work was supported by Grant MH55687 to MJK. We thank Ada Aka, Yuxuan
485 Li, Nicole Kratz, Adam Broitman, Isaac Pedisich, Karl Healey, Patrick Crutchley and Elizabeth Crutchley
486 and other members of the Computational Memory Laboratory at the University of Pennsylvania for their
487 assistance with data collection and preprocessing and Eric Maris, Nora Herweg and Ethan Solomon for
488 helpful comments on a previous version of this manuscript.

489 **Competing Interests** The authors declare that they have no competing financial interests.

490 **Correspondence** Correspondence and requests for materials should be addressed to C.T.W. (email: ctw@cogsci.info).