

Temporal Dynamics of Competition between Statistical Learning and Episodic Memory in Intracranial Recordings of Human Visual Cortex

Brynn E. Sherman,¹ Kathryn N. Graves,¹ David M. Huberdeau,¹  Imran H. Quraishi,²  Eyiymisi C. Damisah,³ and  Nicholas B. Turk-Browne^{1,4}

¹Department of Psychology, Yale University, 2 Hillhouse Avenue, New Haven, CT 06520, ²Department of Neurology, Yale University, 800 Howard Avenue, New Haven, CT 06519, ³Department of Neurosurgery, Yale University, 333 Cedar Street, New Haven, CT 06510, and ⁴Wu Tsai Institute, Yale University, 100 College Street, New Haven, CT 06510

The function of long-term memory is not just to reminisce about the past, but also to make predictions that help us behave appropriately and efficiently in the future. This predictive function of memory provides a new perspective on the classic question from memory research of why we remember some things but not others. If prediction is a key outcome of memory, then the extent to which an item generates a prediction signifies that this information already exists in memory and need not be encoded. We tested this principle using human intracranial EEG as a time-resolved method to quantify prediction in visual cortex during a statistical learning task and link the strength of these predictions to subsequent episodic memory behavior. Epilepsy patients of both sexes viewed rapid streams of scenes, some of which contained regularities that allowed the category of the next scene to be predicted. We verified that statistical learning occurred using neural frequency tagging and measured category prediction with multivariate pattern analysis. Although neural prediction was robust overall, this was driven entirely by predictive items that were subsequently forgotten. Such interference provides a mechanism by which prediction can regulate memory formation to prioritize encoding of information that could help learn new predictive relationships.

Key words: iEEG; learning; memory; prediction

Significance Statement

When faced with a new experience, we are rarely at a loss for what to do. Rather, because many aspects of the world are stable over time, we rely on past experiences to generate expectations that guide behavior. Here we show that these expectations during a new experience come at the expense of memory for that experience. From intracranial recordings of visual cortex, we decoded what humans expected to see next in a series of photographs based on patterns of neural activity. Photographs that generated strong neural expectations were more likely to be forgotten in a later behavioral memory test. Prioritizing the storage of experiences that currently lead to weak expectations could help improve these expectations in future encounters.

Introduction

Long-term memory has a limited capacity; and thus, a major goal of psychology and neuroscience has been to identify factors that determine which memories to store. Well-known factors include attention (Aly and Turk-Browne, 2017), emotion (Dolcos et al., 2017), motivation (Dickerson and Adcock, 2018), stress (Goldfarb, 2019), and sleep (Cowan et al., 2021). Here we further test a novel factor that constrains long-term memory formation: predictive value.

Beyond reliving the past, a key function of memory is that it allows us to predict the future (Schacter et al., 2012). When faced with a new experience, we draw on related experiences from the past to know what is likely to happen when and where (De Brigard, 2014; Biderman et al., 2020). This knowledge is the result of statistical learning, which identifies patterns or regularities in

Received Apr. 7, 2022; revised Oct. 10, 2022; accepted Oct. 13, 2022.

Author contributions: B.E.S. and N.B.T.-B. designed research; B.E.S., K.N.G., D.M.H., I.H.Q., and E.C.D. performed research; B.E.S. analyzed data; B.E.S. wrote the first draft of the paper; B.E.S., K.N.G., D.M.H., I.H.Q., E.C.D., and N.B.T.-B. edited the paper; B.E.S. and N.B.T.-B. wrote the paper.

This work was supported by National Institutes of Health Grant R01 MH069456 to N.B.T.-B.; Canadian Institute for Advanced Research to N.B.T.-B.; and National Science Foundation GRFP Grant to B.E.S. We thank the patients who participated in this study; Kun Wu for providing the electrode reconstructions; Christopher Benjamin for helping to recruit patients and coordinate testing; Richard Aslin and Sami Yousif for helpful conversations; and Gregory McCarthy for advice about data collection and analysis, as well as for feedback on the manuscript.

The authors declare no competing financial interests.

Correspondence should be addressed to Brynn E. Sherman at brynn.sherman@yale.edu or Nicholas B. Turk-Browne at nicholas.turk-browne@yale.edu.

<https://doi.org/10.1523/JNEUROSCI.0708-22.2022>

Copyright © 2022 the authors

the environment that repeat over time (Sherman et al., 2020; Endress and Johnson, 2021) and form the basis of predictions (De Lange et al., 2018). We hypothesize that the availability of these predictions during encoding affects whether a new memory is formed. Namely, if one of the main objectives of long-term memory is to enable prediction, in the service of adaptive behavior, experiences that already generate a prediction may not need to be encoded. In contrast, experiences that yield uncertainty about what will happen next may be more important to store because they can help learn over time what should have been expected. This is distinct from whether an experience being encoded was itself expected or unexpected, which also affects subsequent memory (Greve et al., 2017; Bein et al., 2021); rather, we argue that the process of generating a prediction based on the experience impedes its encoding.

We term this ability of an experience to generate a prediction its predictive value. We previously presented some suggestive evidence for predictive value as an encoding factor. In a statistical learning study with images presented in temporal pairs, subsequent memory for the first item in a pair was impaired relative to unpaired control items (Sherman and Turk-Browne, 2020). Because the first item in a pair was always followed by the second item, it could have enabled a prediction of the second item and thus had predictive value.

However, this prior study was not able to directly link the predictive value of an item during encoding to subsequent memory. From the behavioral data alone (in which prediction was not directly measured), it was unclear whether the memory impairment for the first item originated at the time of encoding or emerged in later stages, such as consolidation or retrieval. For example, the first item might have been encoded well, but when this item was probed in the later memory test, its association with the second item interfered with recognition. Although an fMRI experiment provided some evidence of prediction during encoding — the category of the second item could be decoded during the first — the poor temporal resolution fMRI muddled this interpretation. The decoded neural signals were recorded during or after the second item and shifted backward in time based on assumptions about the hemodynamic lag. Methods with better temporal resolution could provide more precise linking between neural signals and experimental events, allowing for more direct measurement of predictions.

Additionally, in our prior work, we only found a relationship between prediction and encoding across participants. Average fMRI evidence for the category of second items during first items was negatively associated with overall memory performance for first items. However, this could reflect a generic individual difference — that participants who make more predictions tend to have worse memory — rather than prediction having a mechanistic effect on encoding. According to the latter

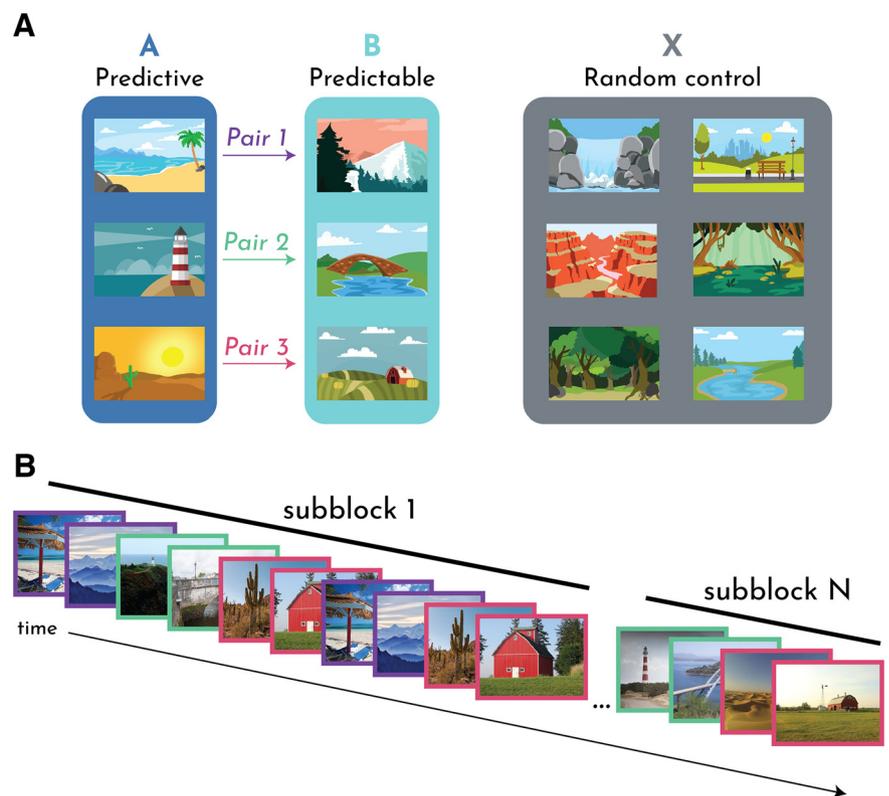


Figure 1. Task design. **A**, Example scene category pairings for 1 participant. Three of 12 categories were assigned to Condition A. Each A category was reliably followed by one of three other categories assigned to Condition B to create pairs. The remaining six categories assigned to Condition X were not paired. Participants viewed the A and B (Structured) and X (Random) categories in separate blocks of the task. **B**, Example stimuli from the Structured block. Participants passively viewed a continuous stream of scenes. Each scene was shown for 267 ms, followed by an ISI of 267 ms with only a fixation cross on the screen. The stream was segmented into subblocks. The same exemplar of each category was presented 4 times per subblock, and new exemplars were introduced for the next subblock. For the Structured block, the category pairs remained consistent across subblocks. Colored frame represents category pairs, corresponding to the A-B pairs (and colored arrows) in subpanel A.

account, whether a participant remembers or forgets a given item should depend on whether that item triggered a prediction during its encoding. This requires testing for a relationship between prediction and encoding across items within participant. Time-resolved methods with denser sampling of individual trials could better enable trial-level estimates of prediction necessary for within-participant subsequent memory analyses.

The present study addresses these issues to better establish predictive value as an encoding factor. We combine intracranial EEG (iEEG) with multivariate pattern analysis, allowing us to measure neural predictions in a time-resolved manner and link them to subsequent behavioral memory across trials. Epilepsy patients viewed a rapid stream of scene photographs across blocks of a statistical learning task. The scenes consisted of unique exemplars from various categories (e.g., beaches, mountains, waterfalls) that differed by block. In the Random blocks, the order of “control” (Condition X) categories from which the exemplars were drawn was random. In the Structured blocks, the categories were paired such that exemplars from “predictive” (Condition A) categories were always followed by exemplars from “predictable” (Condition B) categories (Fig. 1A). Patients were not informed of these conditions or the existence of category pairs, which they learned incidentally through exposure (Brady and Oliva, 2008). The items from each category were presented in subblocks that changed after four presentations (Fig.

Table 1. Patient information^a

ID	Age (yr)	Sex	nElec (vis)	Implant	Data collected	Notes
1	19	F	203 (21)	R G/S/D	2S, 2R	R2 mem data not usable (D)
2	26	F	163 (59)	L G/S/D	2S, 2R	—
3	43	F	172 (10)	Bi D	1S, 2R	—
4	61	F	136 (0)	Bi D	1S, 1R	neural mem data not usable (T)
5	31	M	152 (31)	L G/S/D	2S, 2R	R1 encoding data not usable (T)
6	69	F	92 (7)	L D	2S, 2R	—
7	33	M	232 (22)	Bi D	1S, 1R	—
8	31	F	192 (20)	Bi D	2S, 2R	no mem data collected (C)
9	56	F	192 (36)	Bi D	2S, 2R	R1 encoding data not usable (T)
10	53	M	148 (0)	Bi D	2S, 2R	—

^aDescription of patient participation. nElec (vis), the total number of electrode contacts, followed by the number of visual electrode contacts. Implant: R, right-sided implant; L, left-sided implant; Bi, bilateral implant; G, grid; S, strip; D, depth. Data collected: the number of runs for each condition collected (S, Structured; R, Random). Notes: which runs (if any) were excluded from given analyses and why: D, patient distraction (e.g., a clinician coming in and disrupting testing); T, trigger issue (i.e., an error with the equipment such that we could not align individual trials to our neural signal); C, computer error (e.g., the computer crashed).

1B). After both blocks, patients completed a recognition memory test for the exemplars from the stream.

To track statistical learning in the brain, we used neural frequency tagging (Batterink and Paller, 2017; Choi et al., 2020; Henin et al., 2021). We quantified the phase coherence of oscillations at the frequency of individual items (present in both Random and Structured blocks) and at half of that frequency reflecting groupings of two items (present only in Structured blocks with category pairs). To measure prediction during encoding, we used multivariate pattern similarity (Kok et al., 2014, 2017; Demarchi et al., 2019; Aitken et al., 2020). We first created a template pattern for each scene category based on the neural activity it evoked in visual contacts. We then quantified the expression of these categories during statistical learning, defining prediction as evidence for the second category in a pair evoked by items from the first category.

Although the hippocampus may be the nexus of competition between statistical prediction and episodic encoding (Schapiro et al., 2017; Sherman and Turk-Browne, 2020), hippocampal signals may be relayed and reinstated throughout the cortical hierarchy (Bosch et al., 2014; Tanaka et al., 2014; Hindy et al., 2016; Danker et al., 2017; Aitken and Kok, 2022; Clarke et al., 2022), enabling the robust measurement of learning-related prediction (Kok et al., 2014; Ekman et al., 2017; Kok et al., 2017; Kim et al., 2020) and frequency tagging (Henin et al., 2021) in visual cortex. This allowed us to test our hypotheses robustly in epilepsy patients whose clinical care resulted in extensive electrode coverage in visual cortex but not the hippocampus.

In sum, by assessing iEEG signals during the rapid presentation of scenes, we measured the neural representations underlying statistical learning and prediction, and linked these online learning measures to offline memory, revealing how predictive value constrains memory encoding.

Materials and Methods

Participants

We tested 10 participants (7 female; age range: 19–69 years) who had been surgically implanted with intracranial electrodes for seizure monitoring. Decisions on electrode placement were determined solely by the clinical care team to optimize localization of seizure foci. Participants were recruited through the Yale Comprehensive Epilepsy Center. Participants provided informed consent in a manner approved by the Yale University Human Subjects Committee.

A summary of patient demographics, clinical details, and research participation can be found in Table 1. Given electrode coverage and

usable data, we retained 9 patients in the behavioral analyses, 8 patients in the neural frequency tagging analyses, and 7 patients in the neural category evidence analyses.

iEEG recordings

EEG data were recorded on a NATUS NeuroWorks EEG recording system. Data were collected at a sampling rate of 4096 Hz. Signals were referenced to an electrode chosen by the clinical team to minimize noise in the recording. To synchronize EEG signals with the experimental task, a custom-configured DAQ was used to convert signals from the research computer to 8-bit “triggers” that were inserted into a separate digital channel.

iEEG preprocessing

iEEG preprocessing was conducted in FieldTrip (Oostenveld et al., 2011). A notch filter was applied to remove 60 Hz line noise. No rereferencing was applied, except for 1 patient, whose reference was in visual cortex, resulting in a visual-evoked response in all electrodes; for this patient, we rereferenced the data to a white matter contact in the left anterior cingulate cortex. Data were downsampled to 256 Hz and segmented into trials using the triggers.

Electrode selection

Patients' electrode contact locations were identified using their postoperative CT and MRI scans. Reconstructions were completed in BioImage Suite (Papademetris et al., 2006) and were subsequently registered to the patient's preoperative MRI scan, resulting in contact locations projected into the patient's preoperative space. The resulting files were converted from the Bioimagesuite format (.MGRID) into native space coordinates using FieldTrip functions. The coordinates were then used to create an ROI in FSL (Jenkinson et al., 2012), with the coordinates of each contact occupying one voxel in the mask (Fig. 2).

For purposes of decoding scene categories, we were specifically interested in examining visually responsive contacts (Walther et al., 2009). We defined visual cortex on the MNI T1 2 mm standard brain by combining the Occipital Lobe ROI from the MNI Structural Atlas and the following ROIs from the Harvard-Oxford Cortical Structural Atlas: inferior temporal gyrus (temporo-occipital part), lateral occipital cortex (superior division), lateral occipital cortex (inferior division), intracalcarine cortex, cuneal cortex, parahippocampal gyrus (posterior division), lingual gyrus, temporal occipital fusiform cortex, occipital fusiform gyrus, supracalcarine cortex, and occipital pole. Each ROI was thresholded at 10% and then concatenated together to create a single mask of visual cortex.

To identify which contacts to include in analyses on a per-patient basis, this standard space visual cortex mask was transformed into each participant's native space. We registered each patient's preoperative anatomic scan to the MNI T1 2 mm standard brain template using linear registration (FSL FLIRT) (Jenkinson and Smith, 2001; Jenkinson et al., 2002) with 12 degrees of freedom. This registration was then inverted and used to bring the visual cortex mask into each participant's native space.

In order to ensure that the visual cortex mask captured the anatomic areas we intended, we manually assessed its overlap between the electrodes and made a few manual adjustments to the electrode definition. For example, because of noise in the registrations between postoperative and preoperative space, as well as from preoperative space and standard space, some grid and strip contacts appeared slightly outside of the brain, despite being on the surface of the patient's brain. Thus, contacts such as these were included as “visual” even if they were slightly outside of the bounds of the mask. Additionally, because of the liberal thresholds designed to capture broad visual regions, some portions of the parahippocampal gyrus area contained the hippocampus. Contacts within mask boundaries but clearly in the hippocampus were excluded.

Experimental design

Participants completed the experiment on a MacBook Pro laptop while seated in their hospital bed. The task consisted of up to four runs: two runs of the Structured block and two runs of the Random block. We aimed to collect all four runs from each patient but required a minimum

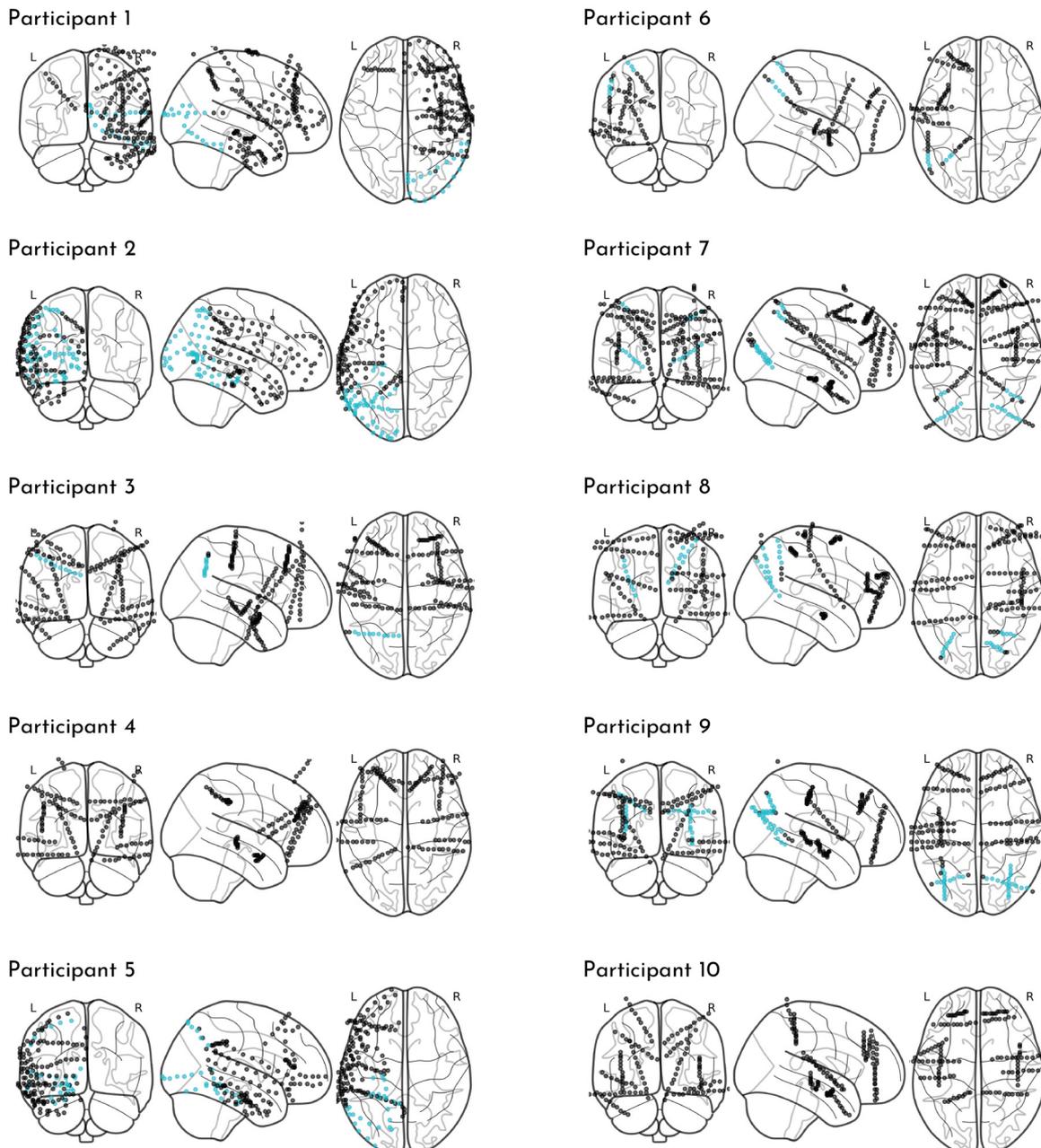


Figure 2. Electrode coverage. The contact locations on the grid, strip, and/or depth electrodes for each participant are plotted as circles in standard brain space. Contacts colored in blue were localized to the visual cortex mask.

of one run per condition for subject inclusion. Given that the order of structured versus random information can impact learning (Jungé et al., 2007; Gebhart et al., 2009), the run order was counterbalanced within and across participants (i.e., some participants received Structured-Random-Random-Structured and others Random-Structured-Structured-Random). Participants completed the runs across 1–3 testing sessions based on the amount of testing time available between clinical care, family visits, and rest times.

Each run consisted of an encoding phase and a memory phase. During the encoding phase, participants viewed a rapid stream of scene images, during which they were asked to passively view the scenes. Participants were told that their memory for the scenes would be tested to encourage them to pay close attention. Each scene was presented for 267 ms, followed by a 267 ms interstimulus interval (ISI) period during which a fixation cross appeared in the center of the screen. These short presentation times were chosen to optimize the task for the frequency

tagging analyses, which involves measuring neural entrainment to stimuli.

Within each run, participants viewed a series of images from a set of six scene categories. There were six categories in the Structured block, and six other categories in the Random block. In the Structured block, the scene categories were paired, such that images from one scene category (A) were always followed by an image from another scene category (B). Thus, A scenes were predictive of the category of the upcoming B scenes, or stated another way, the category of B scenes was predictable given the preceding A scenes. No scene pairs were allowed to repeat back-to-back in the sequence. In the Random block, all six scene categories (X) could be preceded or followed by any other scene category, making them neither predictive nor predictable. No individual scene categories were allowed to repeat back-to-back.

In total, participants viewed 16 exemplars from each category within each run. To assist patients with remembering these briefly presented

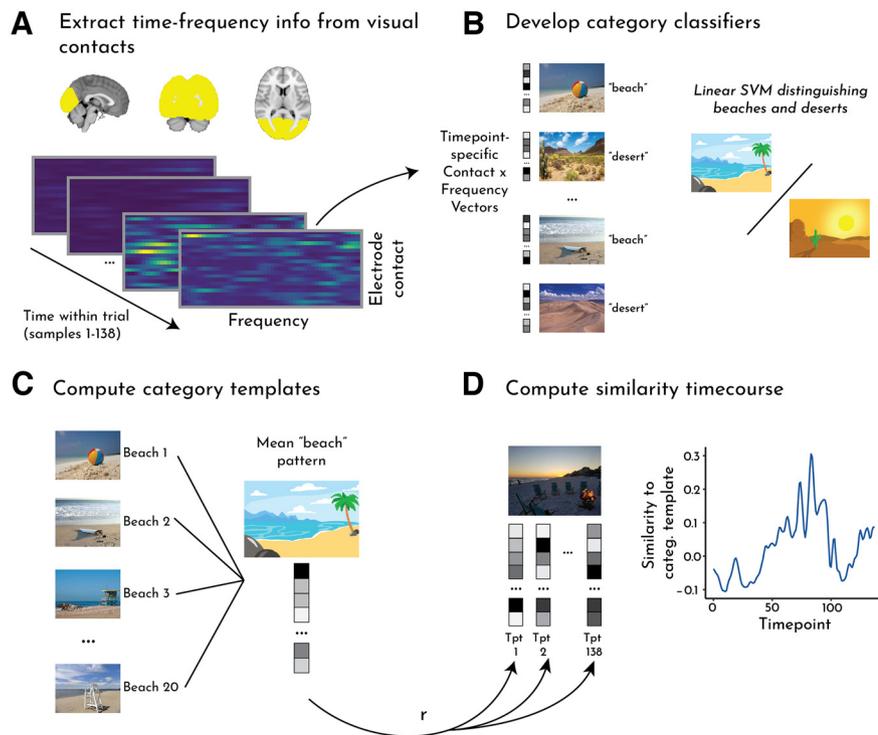


Figure 3. Category evidence analysis pipeline. **A**, A Morlet wavelet approach was used to extract time-frequency information from contacts in visual cortex. This resulted in contact by frequency vectors for every time point of encoding phase and memory phase trials, which served as the neural patterns for subsequent analysis steps. **B**, To identify the neural patterns that distinguished between categories, we ran a series of binary classifiers for every pair of categories from the memory phase trials. These classifiers were trained on the contact by frequency vectors for a single time point or set of time points. The classifiers were then tested on time points from held-out data. **C**, After a series of feature selection steps, we chose the per-participant top-*N* time point set that produced the best classification accuracy, and then averaged contact by frequency vectors across those time points (across all exemplars of a given category) to create a “template” of neural activity for each category. **D**, We then correlated the template for each category from the memory phase with the contact by frequency vector at each time point of each trial/exemplar from that category during the (independent) encoding phase, yielding a time course of pattern similarity reflecting neural category evidence.

images, each individual exemplar was shown 4 times within a run. Thus, each run was comprised of 16 “subblocks” during which the same set of six exemplar images was repeated 4 times (384 trials total). Within each subblock, the order of the pairs/images was randomized, with the constraints described above of no back-to-back repetitions. The individual exemplars changed after each subblock, but the category relations were held constant in the Structured block. Participants were not informed of these category pairings, and thus had to acquire them through exposure.

At the end of each run, participants completed a memory test. Participants were presented with all 96 unique images from the encoding phase, intermixed with 24 novel foils from the same categories (4 foils/category). Participants first had to indicate whether the image was old, meaning it was just presented in that run’s encoding phase, or new, meaning that they had not seen that image at all during the experiment. Following their old/new judgment, participants were asked to indicate their confidence in their response, on a scale of 1 (very unsure) to 4 (very sure). Participants had up to 6 s to make each old/new and confidence judgment. We quantified episodic memory performance using *A'*, a nonparametric measure which takes into account hit rate (HR) and false alarm rate (FA) (Grier, 1971), as follows:

$$A' = .5 + (HR - FA) * (1 + HR - FA) / (4 * HR * (1 - FA))$$

Frequency tagging analyses

We conducted a phase coherence analysis to identify electrode contacts that entrained to image and pair frequencies (Henin et al., 2021). For both Structured and Random blocks, the raw signals were concatenated

across runs (if more than one per block type) and then segmented into subblocks comprising 24 trials with the four repetitions per exemplar. We then converted the raw signals for each subblock into the frequency domain via fast Fourier transform and computed the phase coherence across subblocks for each electrode using the formula $\left[\frac{1}{N} \sum^N (\cos \phi) \right]^2 + \left[\frac{1}{N} \sum^N (\sin \phi) \right]^2$. Notably, by computing phase coherence between subblocks, we collapsed over the contribution of individual exemplars that repeated within subblock. In other words, entrainment in this analysis was driven by phase-locking that generalized across exemplars. Phase coherence was computed separately for each contact in the visual cortex mask, and we then averaged across contacts within participant. We focused on phase coherence at the frequency of image presentation (534 ms/image; 1.87 Hz) and pair presentation (1.07 s/pair; 0.93 Hz).

Category evidence analyses

We used a multivariate pattern similarity approach to assess the time course of category responses. We identified patterns of multivariate activity associated with each category across contacts, frequencies, and time. These category patterns, or “templates,” were defined during the memory phase of the dataset. This was important because the order of categories was random during the memory phase, allowing for an independent assessment of each category across condition regardless of any pairings. We then used these templates to examine category-specific evoked responses during the encoding phase, to assess the presence and timing of category evidence (e.g., for the onscreen category or the upcoming category). The following subsections explain this approach in detail.

Frequency decomposition.

We used a Morlet Wavelet approach to decompose raw signals into time-frequency information (Fig. 3A). We convolved our data with a Complex Morlet Wavelet (cycles = 4) at each of 50 logarithmically spaced frequencies between 2 and 100 Hz to extract the power time course at each of these 50 frequencies. This analysis was done separately for each encoding and memory phase of each run, and the data were z-scored across time within each frequency and contact. This procedure was applied across the unsegmented time courses; we then subsequently carved the time-course into trials using the triggers, yielding a vector of frequency and contact information at each time point within a trial.

Subsequent analyses required that each trial have the same number of time points. However, memory trials were variable lengths, as participants had up to 6 s to respond. There was also slight variability in the encoding trials (most trials were 138 samples long, but some were 136 or 137 samples). To account for this, we considered only the first 138 samples of each memory trial and treated each encoding trial as having 138 samples (interpolating missing time points by averaging the last sample of the trial with the first sample of the next trial).

Category decoding. First, we verified that the multivariate patterns contained category-specific information. We constructed a set of 30 binary classifiers to distinguish among two categories of a given condition (Fig. 3B): A1-A2, A1-A3, A1-B1, A1-B2, A1-B3, A2-A3, A2-B1, A2-B2, A2-B3, A3-B1, A3-B2, A3-B3, B1-B2, B1-B3, B2-B3, X1-X2, X1-X3, X1-X4, X1-X5, X1-X6, X2-X3, X2-X4, X2-X5, X2-X6, X3-X4, X3-X5, X3-X6, X4-X5, X4-X6, X5-X6. We used a linear support vector machine approach using the SVC function in Python’s scikit-learn module, with a penalty parameter of 1.00. We used all of the trials (both old and new exemplars of a category) from the memory runs to train and test the

classifiers and build the subsequent category templates. Thus, there were 20 samples per category for participants who had one run of a condition and 40 samples per category for participants who had two runs of a condition. We split these samples into two-thirds training and one-third test (all within the memory phase), and iterated over the three train-test splits.

First, we independently trained classifiers on a single time point (each of the 138 time points within a trial) and tested each classifier on all 138 time points at test. To validate that we were able to discriminate the categories above chance, we averaged over all train-test combinations and computed overall classification accuracy.

Feature selection. We next aimed to identify the set of time points that produced the best category discrimination. We reasoned that time within a trial would be an important contributor to variance in discriminability, as we would not necessarily expect that time points very early on in a trial (immediately after image onset) would produce high discrimination between categories. We also reasoned that the best time point(s) may differ from participant to participant depending on their electrode coverage. Therefore, we devised a participant-specific time point feature selection process. Importantly, these feature selection steps were conducted within the memory phase data (the same data on which the templates were trained), which were independent of the test data of interest (encoding phase data).

Using the classifier output described above, we averaged the classification over the 138 test time points to assess how well training at every time point generalized to all other time points within a trial. We conducted this analysis for all 30 classifiers and averaged performance over classifiers, yielding a mean classification performance associated with each training time point. For each participant, we then computed the rank order of time points with respect to their classification, such that the first ranked time point was the one that yielded the highest classification, and the last ranked (138th) time point is the one that yielded the lowest classification.

To identify the set of training time points producing the best category classification for a given participant, we repeated the pairwise classification procedure above iteratively training on an increasing number of time points, adding from highest to lowest ranked. Thus, these classifiers ranged from training on the single top time point, to all 138 time points. We again conducted this analysis for all 30 classifiers and averaged performance across them, yielding a mean classification performance associated with the 138 sets of top-*N* time points. We ranked this classification performance again to determine which number of top time points produced the highest classification. This number was used to define the templates.

Template correlations. Using the set of training time points for each participant determined in the feature selection process, we then computed a neural template for each category (Fig. 3C). We extracted the pattern of activity (i.e., a vector containing electrode contact, time, and frequency) for all instances of a given category during the memory phase, including both old and new images. We then averaged over the time points in that participant's training set. The resulting category pattern vector retained spatial (contact) and frequency information.

To assess the time course of neural evidence for a category during the encoding phase, we extracted the pattern of activity (contact and frequency) for each time point of every trial of that category (Fig. 3D). We computed the Pearson correlation between the template and the activity pattern separately for each time point within a trial, yielding a time course of similarity to the template. The resulting Pearson correlation values were Fisher transformed into *z* values.

We were interested in characterizing the time course of a category response not only while that category was on the screen, but also during the surrounding trials. We may observe evidence for a category before it appears, if it can be predicted (as hypothesized for B), or after it disappears, if its representation lingers. Thus, we assessed the time course over a window comprising the onscreen category's trial ("Current") and the two neighboring trials ("Pre" and "Post" trials). To quantify the response, we subtracted a baseline of average evidence for the other categories of the same condition (e.g., for category A1, how much evidence is there for A1 relative to categories A2 and A3?). For the X categories, which could appear in any order, we ensured that the categories included in the baseline did not appear during the "Pre" and "Post" trials. This baselining approach was important for ensuring that effects were not

driven by a generic evoked response (to any category), but rather by specific evidence for the relevant category.

We quantified how template similarity changed over time within trial by splitting the trials into "ON" and "ISI" epochs. The ON epoch refers to the part of the trial when the image was on the screen (the first 69 samples, or 267 ms). The ISI epoch refers to the part of the trial after the image disappeared from the screen during the interstimulus fixation cross (the second 69 samples, or latter 267 ms).

Subsequent memory. To assess how variance in category evidence across trials related to memory outcomes for those trials, we examined predictive and onscreen representations separately for subsequently remembered versus forgotten trials. We conducted this analysis separately for memory of A (as a function of Perceived evidence for A during A and Predicted evidence for B during A) and for memory of B (as a function of Perceived evidence for B during B and Predicted evidence for B during A). Because each image was shown 4 times, we first averaged the Perceived and Predicted evidence over these four trials. We considered the ISI epoch of each trial, as this was the epoch in which we observed reliable evidence for the Predicted category B during A. As a control analysis, we repeated these steps for the X trials from the Random blocks.

Alternative classification approaches for feature selection. The category evidence analyses described above rely on a set of binary classifiers trained to distinguish the categories in a given condition (i.e., all combinations of As and Bs in the Structured condition and Xs in the Random condition). However, this approach may lead to interpretational issues. For example, from a binary classifier trained to distinguish two categories (e.g., A1 vs B1), it is difficult to know whether evidence for one category (e.g., A1) reflects the presence of that category (A1) or the absence of the other category (B1). Thus, we replicated all of the above analyses using two alternative approaches.

First, we trained a 6-way classifier to distinguish among all six categories of a given condition (A1-A2-A3-B1-B2-B3 for Structured and X1-X2-X3-X4-X5-X6 for Random). By including more than two classes, this approach addresses the concern that classification accuracy could be driven by the presence or absence of a given category. Second, we retained the binary classification approach but trained classifiers to only discriminate within the A or B categories. That is, instead of 15 classifiers for A/B combinations, there were 6 classifiers (A1-A2, A1-A3, A2-A3, B1-B2, B1-B3, B2-B3). This approach ensures that classification does not mix evidence for predictive versus predicted categories.

For both of these approaches, we used a linear support vector machine approach using the SVC function in Python's scikit-learn module, with a penalty parameter of 1.00 (same as the primary classification approach). We then repeated the same feature selection steps using these alternative classifiers, and used the output of the top-*N* time point analyses to create new templates.

Statistical analysis

For all analyses (both behavioral and neural), statistical significance was assessed using a random-effects bootstrap resampling approach (Efron and Tibshirani, 1986). For each of 10,000 iterations, we randomly resampled participants with replacement and recomputed the mean across participants, to populate a sampling distribution of the effect. This sampling distribution was used to obtain 95% CIs and perform null hypothesis testing. We calculated the *p* value as the proportion of iterations in which the resampled mean was in the wrong direction (opposite sign) of the true mean; we then multiplied these values by 2 to obtain a two-tailed *p* value. All resampling was done in R (version 3.4.1), and the random number seed was set to 12345 before each resampling test. This approach is designed to assess the reliability of effects across patients: a significant effect indicates that which patients were resampled on any given iteration did not affect the result, and thus that the patients were interchangeable and the effect reliable across the sample.

Results

Memory behavior

We first assessed overall performance in the recognition memory test to verify that participants were able to encode the images

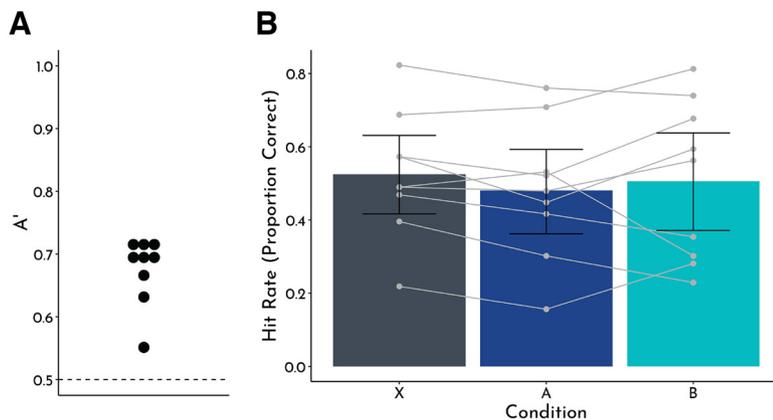


Figure 4. Behavioral results. **A**, Overall memory performance collapsed across conditions. Circle represents A' (a sensitivity measure for recognition memory) for each participant. All participants were above chance (0.5). **B**, Hit rate as a function of condition (A: predictive; B: predictable; X: control). Bars represent group means. Errors bars indicate bootstrapped 95% CI across participants. Individual participant data are overlaid with the gray circles and lines.

into memory. We computed A' , a nonparametric measure of sensitivity (Grier, 1971), from test judgments for items from both Structured and Random blocks. All participants had an A' above the chance level of 0.5 (mean = 0.68; 95% CI = [0.64, 0.70], $p < 0.001$; Fig. 4A) indicating reliable memory. This was driven by a higher hit rate (mean = 0.51) than false alarm rate (mean = 0.32; difference 95% CI = [0.14, 0.23], $p < 0.001$). The proportions of items that were subsequently remembered (hit rate) or forgotten (1-hit rate, or misses) were roughly matched on average, yielding balanced power for within-subject subsequent memory analyses.

We then assessed how statistical learning affected recognition memory. Based on our prior work (Sherman and Turk-Browne, 2020), we hypothesized that the hit rate for items from the predictive A categories in the Structured blocks would be lower than the hit rate for items from the control X categories in the Random blocks. Indeed, we replicated this key behavioral finding (Fig. 4B), with a significantly lower hit rate for A (mean = 0.48) than X (mean = 0.52; difference 95% CI = [−0.076, −0.010], $p = 0.012$). The hit rate for B (mean = 0.51) did not differ from A (difference 95% CI = [−0.10, 0.059], $p = 0.51$) or X (difference 95% CI = [−0.094, 0.053], $p = 0.66$).

The false alarm rate for X (mean = 0.36) was numerically higher than A (mean = 0.28; difference 95% CI = [−0.0023, 0.16], $p = 0.064$); X was significantly higher than B (mean = 0.29; difference 95% CI = [0.0069, 0.13], $p = 0.028$), although A and B did not differ (difference 95% CI = [−0.074, 0.056], $p = 0.82$). Unlike the higher hit rate for X than A, which was specifically hypothesized based on prior work, the marginally higher false alarm rate for X than A was not expected or consistent with previous experiments. Nevertheless, this complicates interpretation of the hit rate difference as impaired memory for A versus X. One difference from the prior study is the blocking of Structured (A,B) and Random (X) categories, which may have allowed for differences in strategy or motivation between conditions. Nevertheless, the main memory hypotheses in the current study rest within the A condition (i.e., which A items are remembered vs forgotten as a function of B prediction), rather than on overall condition-wide differences with X (or B).

We additionally examined the time course of these memory effects by sorting the items into subblocks. If the deficit in memory for A items arises from the predictive value that they confer,

we might expect that this effect will emerge over time as learning occurs (Sherman and Turk-Browne, 2020). We focused this analysis on the first Structured run of the encoding phase for each participant, to equate the amount of data and corresponding opportunity for learning across participants (some had one run, others two). We quantified change over time for each participant as the Spearman rank correlation of subblock number with hit rate for A (averaged across items in each subblock), expecting a negative correlation. The resulting within-participant relationship was not reliable at the group level (mean $\rho = -0.038$; 95% CI = [−0.27, 0.19], $p = 0.77$). This null effect of a learning trajectory stands in contrast with what we observed in Sherman and Turk-Browne (2020), perhaps related to the smaller number of participants or differences in task design (e.g., the use of subblocks) in the current study.

Neural frequency tagging

To provide a neural check of statistical learning of the category pairs in the Structured blocks, we measured entrainment of neural oscillations in visual electrode contacts to the frequency of individual images and image pairs (Fig. 5A). We expected strong entrainment at the image frequency in both the Structured and Random blocks, as this reflects the periodicity of the sensory stimulation. Critically, we hypothesized that there would be greater entrainment at the pair frequency in Structured compared with Random blocks. This provides a measure of statistical learning because the pairs only exist when participants extract regularities over time in the transition probabilities between categories in the Structured blocks.

Consistent with our hypotheses and prior work (Henin et al., 2021), there were distinct peaks in phase coherence at both the image and pair frequencies in Structured blocks, but only at the image frequency in Random blocks (Fig. 5B). To confirm the reliability of these peaks, we contrasted the coherence at the frequency of interest (image: 1.87 Hz; pair: 0.93 Hz) against a baseline of the coherence at frequencies neighboring each of the frequencies of interest (± 0.078 Hz). At the image frequency, there were reliable peaks in both the Structured (mean difference = 0.46; 95% CI = [0.37, 0.55], $p < 0.001$) and Random blocks (mean difference = 0.42; 95% CI = [0.28, 0.52], $p < 0.001$). At the pair frequency, there was a reliable peak in Structured blocks (mean difference = 0.059; 95% CI = [0.035, 0.084], $p < 0.001$), but not Random blocks (mean difference = −0.0027; 95% CI = [−0.016, 0.0085], $p = 0.68$).

Further, the peak in coherence at the pair frequency in Structured blocks was reliably higher than that in Random blocks (mean difference = 0.058; 95% CI = [0.035, 0.083], $p < 0.001$), confirming that the pair frequency effect was specific to when there was structure in the sequence. There were no differences in coherence at the image frequency across conditions (mean difference = 0.018; 95% CI = [−0.010, 0.048], $p = 0.25$). Together, these results provide strong evidence that visual regions represented the paired categories during statistical learning.

To measure the emergence of these entrainment effects over time, we computed the coherence over an iteratively increasing number of subblocks (Henin et al., 2021). Specifically, we first computed the coherence across the first two subblocks, then the

first three, and so on, up to all 16 subblocks. As in the behavioral time course analyses, we only included the first 16 subblocks per participant (corresponding to the first run of a given condition) to equate the opportunity for learning effects across participants. To quantify neural entrainment, we computed the difference in coherence between the frequency of interest and the two neighboring frequencies (as we did above to establish whether peaks were reliable). We then assessed the reliability of that difference, relative to 0, across participants. We hypothesized that coherence at the pair frequency would emerge over time in the Structured condition, but that coherence at the image frequency would be consistently high, even at early time points.

In the Structured condition, the pair frequency was consistently reliable by the 13th subblock (mean ITC difference = 0.035; 95% CI = [0.0011, 0.071], $p = 0.043$), with each subsequent subblock also exhibiting a reliable peak in coherence at the pair frequency (p values < 0.001; Fig. 5C, left). Confirming that this effect was specific to the Structured condition, we did not find reliable peaks in coherence at the pair frequency across any number of subblocks in the Random condition (p values > 0.30).

In contrast to the pair frequency that required learning, the image frequency should be driven by the stimuli and thus present early in both conditions. Indeed, coherence at the image frequency was reliably high across all numbers of subblocks, in both the Structured and Random conditions (all p values < 0.001; Fig. 5C, right). This lends credence to the interpretation of increasing coherence at the pair frequency over time as reflecting a trajectory of learning.

Given our interpretation that entrainment to the pair frequency reflects statistical learning, and given that we expect our key behavioral effect (impaired memory for predictive A items) to depend on statistical learning, we next asked whether these two effects are related. We calculated this relationship within-participant given the small sample for estimating across-participant relationships. Coherence is necessarily measured across trials; and thus, we could not relate entrainment on a given trial to memory for that trial. Instead, we computed coherence across neighboring subblocks and estimated neural entrainment to the pairs as the difference in coherence at the pair frequency from the two adjacent frequencies. We then related this neural measure to average A hit rate within the latter of the two neighboring subblocks, expecting a negative relationship (stronger pair entrainment associated with worse A memory). For example, the coherence between Subblocks 1 and 2 was used to predict behavioral memory in Subblock 2 (memory in Subblock 1 was excluded from this analysis). The within-participant relationship between neural entrainment to pairs and A memory showed a trend at the group level (mean $\rho = -0.13$; 95% CI = [-0.25, 0.020], $p = 0.089$), although importantly 6 of 7 participants showed a negative

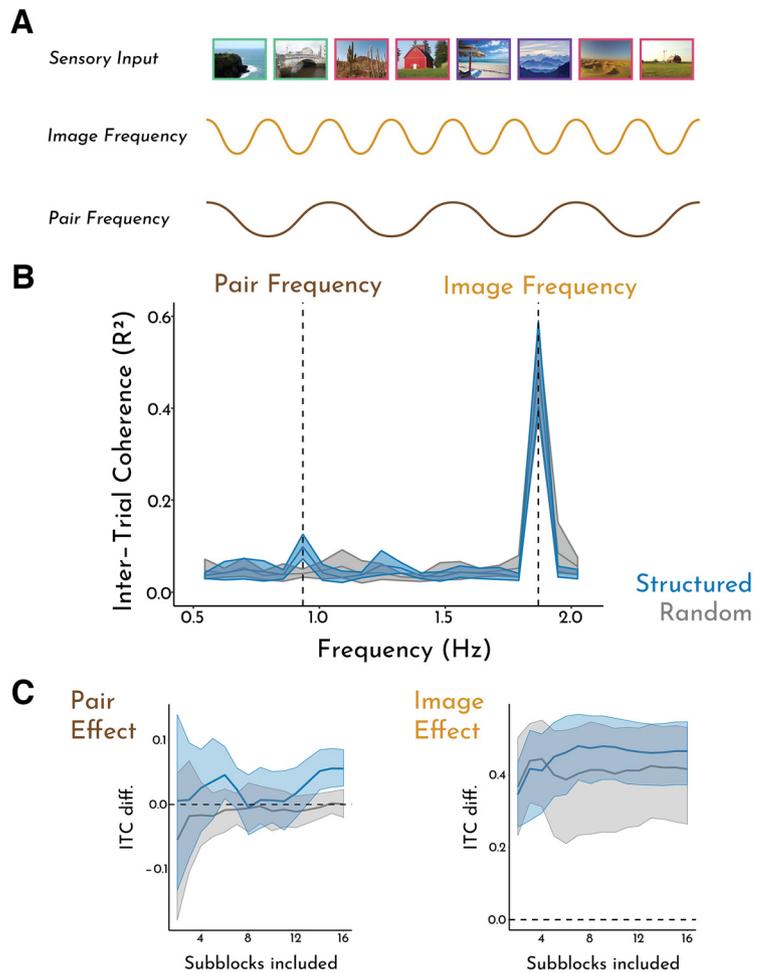


Figure 5. Neural frequency tagging analysis. **A**, Schematic of analysis and hypothesized neural oscillations. We expect entrainment of visual contacts at the frequency of images in both blocks. In the Structured block, we also expect entrainment at the frequency of category pairs. **B**, These hypotheses were confirmed, with reliable peaks in coherence at the image and pair frequencies in Structured blocks but only at the image frequency in Random blocks. **C**, We examined the emergence of entrainment over time by measuring the difference in coherence at the frequency of interest, relative to the two neighboring frequencies, as we iteratively increased the number of subblocks from the start of the run included in the analysis. Left, Coherence at the pair frequency emerged over time in the Structured block (reaching significance by the 13th subblock and beyond) but not in the Random block. Right, Coherence at the image frequency was high in both blocks, regardless of how many subblocks were included. Error bands indicate the 95% bootstrapped CIs across participants.

correlation. We repeated this analysis for the image frequency as a control, and found no relationship between neural entrainment to images and A memory (mean $\rho = -0.072$; 95% CI = [-0.24, 0.087], $p = 0.42$).

Scene category decoding and template creation

The neural frequency tagging for pairs in Structured blocks indicates statistical learning of the pairs. This learning should create predictive value for the items from the A categories, which afford a prediction of the associated B category. To test for these predictive representations, we used a multivariate pattern similarity approach that extracted neural evidence for visual categories. For each category, we created a neural template based on the pattern of time-frequency information evoked by each category across visual contacts. These templates were optimized through a series of steps (described below) for each participant to ensure maximum category discriminability.

First, to verify that the scene categories were indeed discriminable, we developed a series of binary classifiers to distinguish

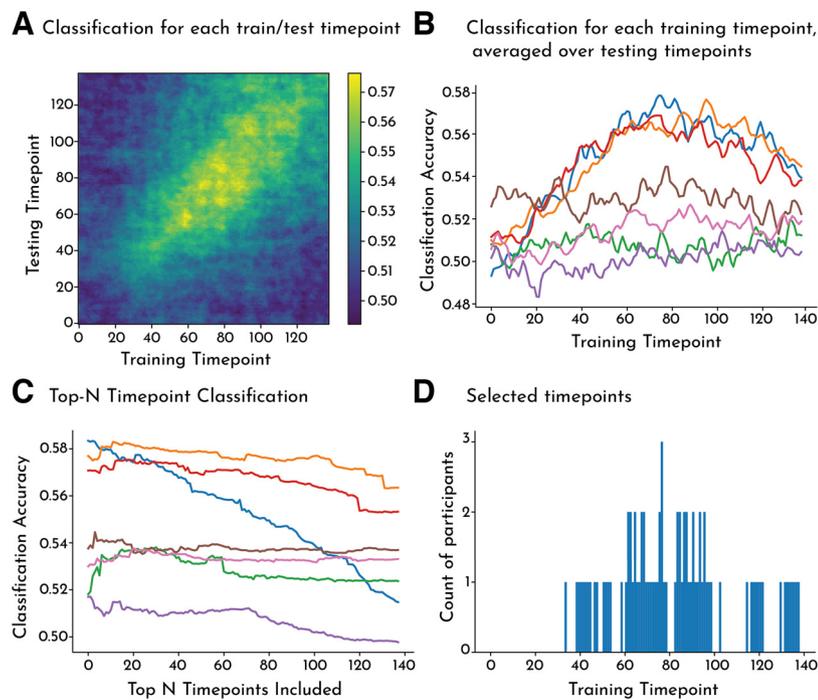


Figure 6. Category decoding and feature selection. **A**, To establish overall category decoding accuracy, we trained and tested binary category classifiers separately for all individual time points, yielding a temporal generalization matrix. **B**, As a first feature-selection step, we computed the average classification accuracy (across pairwise classifiers) for each training time point and participant (colored lines). We then ranked the time points by classification accuracy. **C**, To select the set of time points that produced the best classification for a given participant, we trained and tested the category classifiers on an increasing number of time points, starting with the best-performing time point identified in **B** and iteratively adding time points by rank. We then computed the per-participant average classification accuracy for each set of time points. **D**, Histogram depicting which training time points were selected for template creation for all participants (e.g., count = 3 indicates that that time point was included for 3 of the 7 participants).

among the scene categories. Because we were interested in ultimately selecting the time points that produced the best category discrimination, we trained classifiers on a single time point (each of the 138 time points within a trial) and tested each classifier on all 138 time points at test. Figure 6A illustrates the classification performance across all of these binary classifiers, averaged across participants. At the group level (averaging across all train-test combinations), classification performance was above chance (mean = 0.528; 95% CI = [0.514, 0.542], $p < 0.001$), with each individual participant exhibiting classification performance greater than the chance level of 0.5.

We next aimed to select the training time points (per participant) that exhibited the best category discrimination. For each participant and training time point, we averaged classification accuracy across all test time points (Fig. 6B). We then ranked the training time points by classification accuracy. Next, to find the set of training time points that produced the best classification, we reran our classification procedure, but training on an increasing number of time points, starting with the best-performing time point, and iteratively adding time points per rank. We then computed the per-participant average classification accuracy for each set of time points (Fig. 6C). Verifying that this feature selection approach worked to optimize category discriminability, we indeed found that using the per-participant top- N time points yielded higher classification accuracy than averaging across all time points (mean accuracy = 0.554; 95% CI = [0.536, 0.571], $p < 0.001$); this was independently true for each participant.

We used these per-participant top- N time points to create templates of each category. Figure 6D illustrates the training time points which were included in the templates, for 1 or more participants. To construct the templates, we averaged the contact \times frequency vectors across the top- N time points for all exemplars of a given category. We then aimed to quantify the expression of these category templates during learning (e.g., during the presentation of a predictive A item, is there a representation of the upcoming B item?). However, given that these templates were created from the memory phase, after learning had already occurred, it is important to ensure that the templates of paired categories themselves were not correlated with each other; if so, any effects of prediction during learning could be confounded. At the group level, the templates of paired categories (e.g., A1-B1) were no more correlated than the templates of unpaired Structured categories (e.g., A1-B2; mean difference = 0.024; 95% CI = [-0.019, 0.069], $p = 0.30$) or Random categories (e.g., X1-X2; mean difference = 0.047; 95% CI = [-0.032, 0.127], $p = 0.25$).

Category evidence during learning

To test for evidence of predictive value, we quantified the expression of these templates in the Structured and Random blocks. As a check, we expected clear neural evidence for the category of the item being presented on the screen. Critically, we hypothesized that neural evidence for the upcoming B category would manifest before its appearance, in response to an A exemplar.

We measured these temporal dynamics of neural category evidence by creating a window of three trials centered on the current item: the trial preceding a trial in which the item appeared (“Pre”), the trial during which the item was on the screen (“Current”), and the trial succeeding the trial in which the item appeared (“Post”). For example, if category Pair 1 involved beaches (A1) being followed by mountains (B1), neural evidence for the mountain category was calculated in response to beach exemplars (Pre), mountain exemplars (Current), and exemplars from the categories that could appear next in the Structured sequence (A2 or A3 categories). These evidence values were averaged across the categories from the same condition (e.g., B1, B2, and B3 for Condition B) and plotted over time (Fig. 7A). For statistical analysis, we averaged the neural category evidence for each category across the time points within 6 epochs: when Pre, Current, and Post images were on the screen (ON) and during the fixation period between these trials (ISI; Fig. 7B). We anticipated the evoked response to each image would span ON and ISI periods (as neural processing of the image would take longer than 267 ms), but subdividing in this way allowed us to test for the emergence of predictive evidence of B during the ISI immediately before its onset.

For Current trials (i.e., the trial when the target category was on screen), we found robust (perceptual) evidence for both A and B across both the ON epoch (A: mean = 0.0088; 95% CI = [0.0046, 0.013], $p < 0.001$; B: mean = 0.012; 95% CI = [0.0066, 0.018], $p < 0.001$) and ISI epoch (A: mean = 0.012; 95% CI = [0.0084, 0.015], $p < 0.001$; B: mean = 0.014; 95% CI = [0.0083,

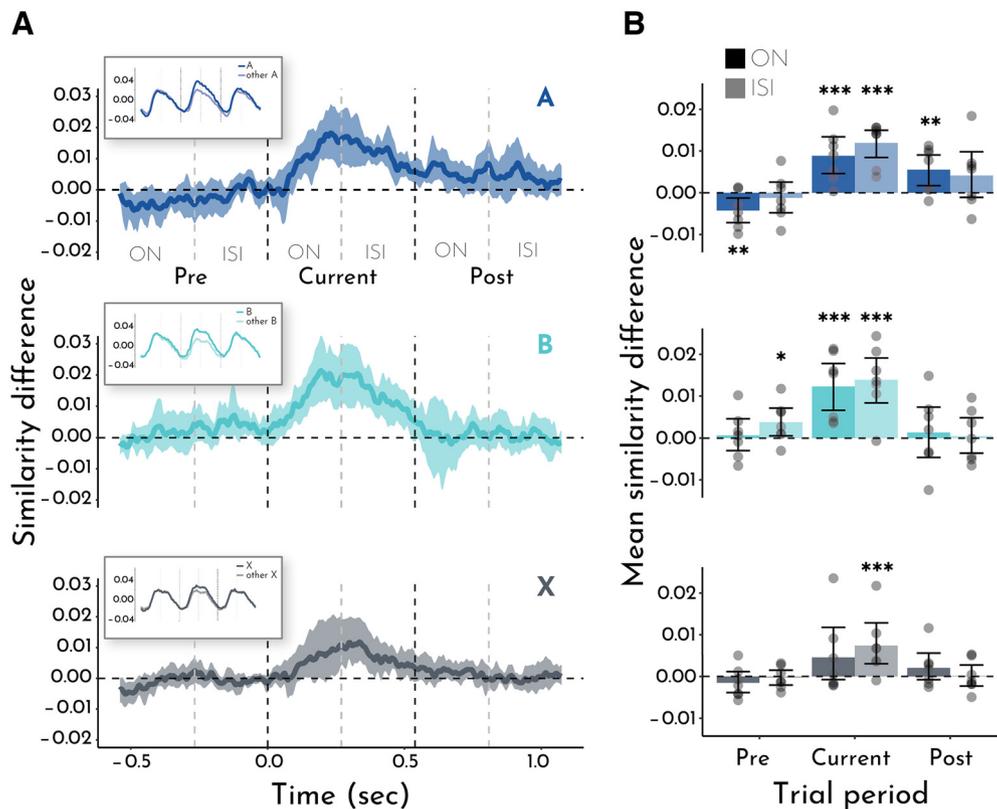


Figure 7. Neural category evidence. **A**, Time course of similarity between patterns of neural activity in visual contacts evoked by exemplars from A (predictive), B (predictable), and X (control) categories and category template patterns for A, B, and X, respectively, baselined to average evidence for the other categories of the same condition. Inset, Raw pattern similarity before baseline subtraction for the category template of interest (dark) and the average of the other category templates from the same condition (light). Error bands were removed for ease of visualization. Current, the trial when the item was presented; Pre, the trial before the item was presented; Post, the trial after the item was presented. For each row/condition, the Pre, Current, and Post trials are compared with the same category template (Current). Error bands represent the bootstrapped 95% CIs across participants (i.e., any time point whose band excludes 0, $p < 0.05$). **B**, Average pattern similarity collapsed across time points within ON (stimulus on screen) and ISI (fixation between stimuli) epochs. Each dot represents an individual participant. Bars represent the means across participants. Error bars indicate the bootstrapped 95% CIs. * $p < 0.05$. ** $p < 0.01$. *** $p < 0.001$.

0.019], $p < 0.001$). Neural evidence for X categories from Random blocks was not reliable during the ON epoch (mean = 0.0046, 95% CI = [−0.00075, 0.012], $p = 0.13$) but became robust later in the trial during the ISI epoch (mean = 0.0074; 95% CI = [0.0030, 0.013], $p < 0.001$). There was greater evidence for B than X categories during both ON (mean difference = 0.0077; 95% CI = [0.00058, 0.015], $p = 0.031$) and ISI epochs (mean difference = 0.0065; 95% CI = [0.00061, 0.012], $p = 0.031$). Considering X as a baseline, this difference shows enhanced perceptual processing of predictable categories. Neural evidence did not differ between A and B categories (p values > 0.38) or A and X categories (p values > 0.28).

For Pre trials (i.e., the trial before the target category appeared), we found the hypothesized predictive neural evidence for the B categories during the ISI epoch (just after its paired A category appeared; mean = 0.0037; 95% CI = [0.00054, 0.0071], $p = 0.019$). B evidence was not present during the ON epoch earlier in the Pre trials (while its paired A category was on screen; mean = 0.00063; 95% CI = [−0.0030, 0.0046], $p = 0.78$); this may reflect the time needed for associative reactivation of the B category after perceptual processing of the A item, or anticipation of the timing when B will appear (at the end of the Pre trials). Further supporting our interpretation that Pre evidence of the B categories reflects prediction, no such evidence was observed for X during ON (mean = −0.0015; 95% CI = [−0.0039, 0.0012], $p = 0.26$) or

ISI epochs (mean = −0.00031; 95% CI = [−0.0021, 0.0015], $p = 0.73$) or for A during the ISI epoch (mean = −0.0012; 95% CI = [−0.0048, 0.0025], $p = 0.53$). There was negative evidence for the upcoming A category during the ON epoch of the Pre trial (mean = −0.0043; 95% CI = [−0.0072, −0.0013], $p = 0.0052$), but this may have been artifactual (see below). When contrasting prediction-related signals across conditions, Pre neural evidence for the B categories during the ISI epoch was reliably greater than X categories (mean difference = 0.0040; 95% CI = [0.00016, 0.0075], $p = 0.042$) and marginally greater than A categories (mean difference = 0.0049; 95% CI = [−0.00051, 0.010], $p = 0.075$).

For Post trials (i.e., the trial after the target category appeared), we found reliable neural evidence for the A categories during the ON epoch (i.e., while its paired B category was on screen; mean = 0.0055; 95% CI = [0.0017, 0.0091], $p = 0.0018$); this effect was not significant during the ISI epoch (mean = 0.0041; 95% CI = [−0.0011, 0.0098], $p = 0.13$). We did not find Post evidence of B or X categories during either ON or ISI epochs (p values > 0.80), nor was Post evidence for A reliably stronger than B or X (p values > 0.16). Positive evidence of A during the Post trial may be related to the negative evidence of A during the Pre trial noted above. Because no back-to-back pair repetitions were allowed, in an A1-B1-A2-B2 trial sequence, A1 and A2 were different categories. A1 evidence during B1 was considered a Post trial for the A condition, whereas A2 evidence during B1

was considered a Pre trial for the A condition. Because A1 was one of two baseline categories for A2 (along with the third A category, A3), Post evidence for A1 during B1 would have been subtracted from Pre evidence for A2, leading to a negative effect. We tested this by comparing evidence for A2 (Pre) and A1 (Post) during B1 to the neutral A3 only. This weakened the negative Pre evidence for A, during ON (mean = -0.0027 ; 95% CI = $[-0.0054, 0.00]$, $p = 0.058$) and ISI epochs (mean = 0.00048 ; 95% CI = $[-0.0022, 0.0038]$, $p = 0.82$). However, the positive Post evidence for A during the ON epoch remained significant (mean = 0.0081 ; 95% CI = $[0.0036, 0.014]$, $p < 0.001$).

The findings above rely on category templates optimized based on a set of binary category classifiers. To ensure that our results are robust to these specific feature selection steps, we reran our analyses using two different approaches for template creation.

First, we created category templates from a 6-way classifier that simultaneously learned to distinguish the patterns from all categories of a condition. As a check, we first confirmed that this method produced the same results for Current items. Indeed, as above, we found reliable evidence for both A and B items, during the ON (A: mean = 0.0095 ; 95% CI = $[0.0056, 0.014]$, $p < 0.001$; B: mean = 0.015 ; 95% CI = $[0.010, 0.019]$, $p < 0.001$) and ISI periods (A: mean = 0.010 ; 95% CI = $[0.0060, 0.014]$, $p < 0.001$; B: mean = 0.014 ; 95% CI = $[0.0085, 0.019]$, $p < 0.001$); evidence for X was reliable during the ISI (mean = 0.0059 ; 95% CI = $[0.0026, 0.0099]$, $p < 0.001$), but not ON periods (mean = 0.0037 ; 95% CI = $[-0.0026, 0.012]$, $p = 0.32$). Critically, we replicated our key finding of predictive B evidence during the Pre-ISI period (i.e., just after its paired A category appeared; mean = 0.0035 ; 95% CI = $[0.00042, 0.0066]$, $p = 0.025$), as well as of lingering A evidence during the Post-ON period (i.e., while its paired B category was on screen; mean = 0.0049 ; 95% CI = $[0.000059, 0.0095]$, $p = 0.049$).

Second, we retained the binary classification approach but limited the classifiers to category comparisons within A or within B, such that the classifiers did not learn to discriminate A versus B. Although we expected that this approach would reduce the quality of feature selection by optimizing for fewer category distinctions, it eliminated the possibility that mixing predictive and predicted categories may artificially inflate classification performance. This approach again produced qualitatively similar results, though slightly weaker. We found reliable evidence for both A and B Current items, during the ON (A: mean = 0.0093 ; 95% CI = $[0.0060, 0.013]$, $p < 0.001$; B: mean = 0.013 ; 95% CI = $[0.0076, 0.018]$, $p < 0.001$) and ISI periods (A: mean = 0.010 ; 95% CI = $[0.0063, 0.013]$, $p < 0.001$; B: mean = 0.015 ; 95% CI = $[0.0097, 0.020]$, $p < 0.001$); evidence for X was reliable during the ISI (mean = 0.0078 ; 95% CI = $[0.0045, 0.012]$, $p < 0.001$), but not ON periods (mean = 0.0046 ; 95% CI = $[-0.0012, 0.012]$, $p = 0.17$). Further, we numerically replicated our key finding of predictive B evidence during the Pre-ISI period (mean = 0.0038 ; 95% CI = $[0.00, 0.0080]$, $p = 0.050$), though lingering A evidence during the Post-ON period was no longer reliable (mean = 0.0022 ; 95% CI = $[-0.0034, 0.0081]$, $p = 0.47$).

Together, these results show that statistical learning of the category pairs in Structured blocks affected neural representations in the task. Not only did visual contacts represent the category of the first and second items in a pair while they were being perceived (A and B evidence during ON and ISI epochs of A and B, respectively), but also the first category during the second (A evidence during ON epoch of B) and the second category during

the first (B evidence during ISI epoch after A). This latter effect indicates that the first item in a pair (from A category) had predictive value on average.

We again examined whether these predictive effects emerged over time, in the first run of the Structured condition. For each participant, we computed the Spearman rank correlation of subblock number with the mean predictive evidence for B (averaged across all A items in each subblock), expecting a positive correlation. The resulting within-participant relationship was not reliable at the group level (mean $\rho = 0.012$; 95% CI = $[-0.24, 0.24]$, $p = 0.92$). We also tested for a positive relationship across subblocks between prediction of B during A and neural entrainment for pairs, given that we expect both measures to depend on statistical learning. However, this within-participant relationship was not reliable at the group level (mean $\rho = 0.038$; 95% CI = $[-0.12, 0.19]$, $p = 0.67$), nor was it reliable for neural entrainment to images (mean $\rho = -0.11$; 95% CI = $[-0.29, 0.079]$, $p = 0.25$).

Although we did not observe a clear learning trajectory, we can still leverage variability in prediction across trials to understand the relationship between predictive value and memory.

Subsequent memory analysis

We theorized that items with predictive value are a lower priority for new encoding into episodic memory. Here we test this relationship by comparing neural category evidence for remembered versus forgotten items within participants. That is, although A items had reliable predictive value on average, variability across items may relate to subsequent memory. To the extent that prediction interferes with encoding, we hypothesized that subsequently forgotten A items would elicit evidence for the upcoming B category during their encoding. Critically, in contrast to prior analyses relating entrainment to memory or prediction, which required measurements at the subblock level, here we are able to probe the relationship between prediction and memory at the level of individual trials.

Consistent with our hypothesis, B evidence during the ISI epoch after A (i.e., Predicted category) was negatively related to subsequent A memory (Fig. 8A): forgotten A items yielded reliable B evidence (mean = 0.0092 ; 95% CI = $[0.0023, 0.017]$, $p = 0.0030$), whereas remembered A items did not (mean = 0.0017 ; 95% CI = $[-0.0016, 0.0049]$, $p = 0.31$). In contrast, A evidence during the ISI epoch after A (i.e., Perceived category) was reliable for both remembered (mean = 0.012 ; 95% CI = $[0.0091, 0.015]$, $p < 0.001$) and forgotten (mean = 0.014 ; 95% CI = $[0.0077, 0.021]$, $p < 0.001$) A items. This differential effect of subsequent memory on neural evidence for Perceived versus Predicted categories during the ISI after A was reflected in a significant 2 (evidence category: A, B) by 2 (subsequent memory: remembered, forgotten) interaction ($p < 0.001$). This interaction was driven by a marginal difference in neural evidence for the Predicted B category during encoding of subsequently forgotten versus remembered A items (mean difference = 0.0075 ; 95% CI = $[-0.00046, 0.016]$, $p = 0.065$), but no reliable difference in neural evidence for the Perceived A category by subsequent memory (mean difference = 0.0022 ; 95% CI = $[-0.0050, 0.0094]$, $p = 0.57$).

As a control analysis, we performed the key steps above in the Random blocks. These blocks did not contain pairs, and so we dummy-coded pairs of X items ($X_1 - X_2$ instead of A-B). In contrast to Structured blocks, we did not expect that neural evidence of the “Predicted” X_2 category during the X_1 ISI would relate to subsequent memory for X_1 . Indeed, there was no reliable evidence for the X_2 category for either remembered (mean =

−0.0029; 95% CI = [−0.0069, 0.00084], $p = 0.14$) or forgotten (mean = 0.0011; 95% CI = [−0.0027, 0.0054], $p = 0.57$) X_1 items. In contrast, neural evidence for the Perceived X_1 category during the X_1 ISI was reliable for both remembered X_1 items (mean = 0.010; 95% CI = [0.0039, 0.019], $p < 0.001$) and forgotten X_1 items (mean = 0.0065; 95% CI = [0.0022, 0.012], $p < 0.001$).

We so far focused on the effects of prediction for memory of the item generating the prediction (A), but what is the mnemonic fate of the item being predicted (B), which in this task with deterministic pairs always appeared as expected? Whereas neural category evidence for B during the A ISI (Predicted) was negatively related to subsequent memory for A items, the opposite was true for memory of B items (Fig. 8B): remembered B items were associated with reliable prediction of B (mean = 0.0082; 95% CI = [0.0036, 0.012], $p < 0.001$), but forgotten B items were not (mean = −0.0028; 95% CI = [−0.011, 0.0041], $p = 0.49$). In contrast, and similar to A memory, evidence for B during the B ISI (Perceived) was reliable for both remembered (mean = 0.013; 95% CI = [0.0082, 0.018], $p < 0.001$) and forgotten (mean = 0.014; 95% CI = [0.00096, 0.026], $p = 0.034$) B items. We did not find an interaction between category and memory ($p = 0.22$). However, there was a reliable difference in Predicted B evidence for remembered versus forgotten B items (mean difference = 0.011; 95% CI = [0.00060, 0.021], $p = 0.039$); Perceived B evidence did not differ as a function of memory (mean difference = 0.00064; 95% CI = [−0.014, 0.016], $p = 0.89$).

We repeated the same control analysis of Random blocks, but now focused on subsequent memory for X_2 items (equivalent to B, rather than X_1 memory for A). Neural evidence for the “Predicted” X_2 category during the ISI after X_1 was not reliable for either remembered (mean = 0.0013; 95% CI = [−0.0020, 0.0043], $p = 0.44$) or forgotten (mean = −0.00048; 95% CI = [−0.0030, 0.0017], $p = 0.75$) X_2 items.

We again tested whether our key results generalized to templates created from two alternative classification approaches. Using a 6-way classifier, we replicated the finding that forgotten A items were associated with reliable predictive evidence of B (mean = 0.0075; 95% CI = [0.0015, 0.014], $p = 0.009$), whereas remembered A items were not (mean = 0.0026; 95% CI = [−0.00010, 0.0054], $p = 0.061$). In contrast, forgotten B items were not associated with reliable predictive evidence of B (mean = −0.0046; 95% CI = [−0.016, 0.0037], $p = 0.40$), whereas remembered B items were (mean = 0.0082; 95% CI = [0.0021, 0.015], $p = 0.003$). Using binary classifiers trained to discriminate within A or B categories, we again found that forgotten (mean = 0.0075; 95% CI = [0.00087, 0.016], $p = 0.014$), but not remembered A items (mean = 0.0027; 95% CI = [−0.00086, 0.0061],

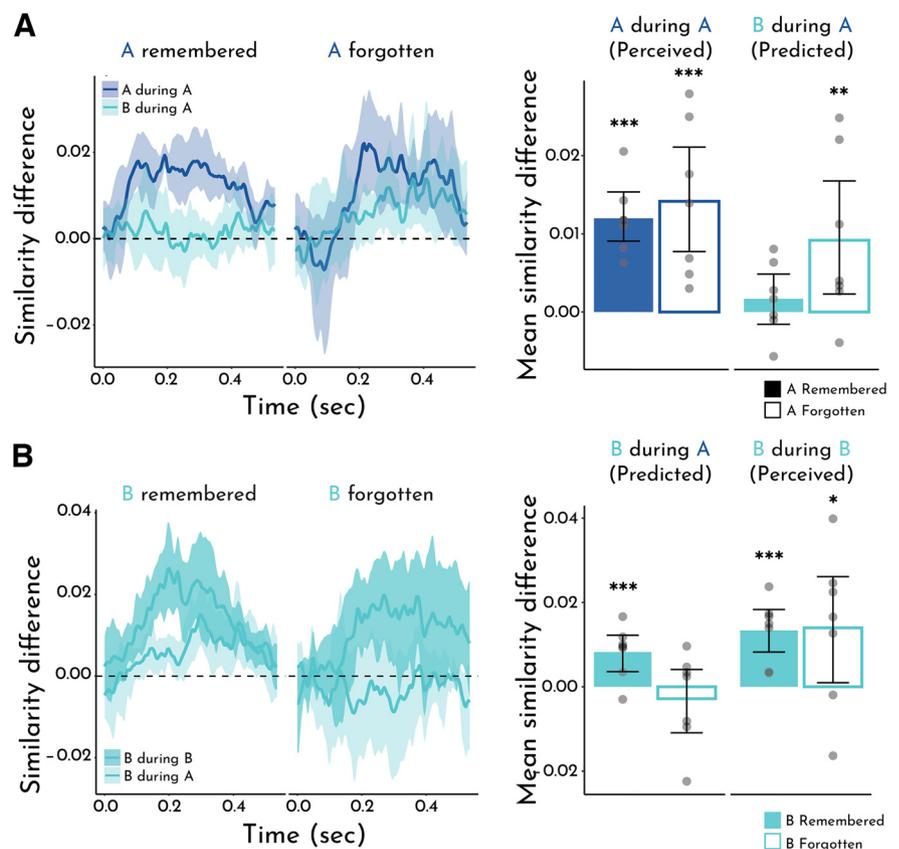


Figure 8. Subsequent memory analysis. **A**, Left, Time course of pattern similarity in visual contacts between the encoding of A items and the category templates for A (Perceived, A during A) and B (Predicted, B during A), as a function of whether the A items were subsequently remembered or forgotten. Right, Pattern similarity averaged within the ISI period, the epoch in which we observed overall evidence of prediction, as a function of subsequent memory for A items (filled bars represent remembered; empty bars represent forgotten). **B**, Left, Time course of pattern similarity in visual contacts between the category template for B and the encoding of A items (Predicted, B during A) and B items (Perceived, B during B), as a function of whether the B items were subsequently remembered or forgotten. Right, Pattern similarity averaged within the ISI period, as a function of subsequent memory for B items. Error shading/bars represent the bootstrapped 95% CI across participants. Each dot represents an individual participant. * $p < 0.05$. ** $p < 0.01$. *** $p < 0.001$.

$p = 0.13$) were associated with reliable predictive evidence of B, and that remembered (mean = 0.0084; 95% CI = [0.0033, 0.013], $p = 0.0016$), but not forgotten B items (mean = −0.0044; 95% CI = [−0.017, 0.0048], $p = 0.47$) were associated with reliable predictive evidence of B.

Together, these results highlight the opposing influence of predictive value on memory for predictive versus predicted items. Namely, prediction of B (during A) is associated with worse memory for predictive A items (suggesting interference between the generation of a prediction and encoding of the current item) but better memory for predicted B items (suggesting that this prediction may potentiate encoding of an upcoming item).

Discussion

This study demonstrates a trade-off between how well an item is encoded into episodic memory and how strong of a future prediction it generates based on statistical learning. We first used frequency tagging to provide neural verification of statistical learning. During a sequence of scene photographs, electrodes in visual cortex represented pairs of scene categories that reliably followed each other, synchronizing not only to the individual scenes but also to the boundaries between pairs. Next, we used

multivariate pattern analysis to assess how the paired categories were represented over time. Items from the first category in a pair elicited a representation of the second category, which grew in strength in advance of the onset of items from the second category. We refer to the ability of an item to generate this predictive representation as its “predictive value.” Critically, by relating these representational dynamics to subsequent memory behavior, we found that forgotten items from the first category triggered reliable predictions during encoding whereas remembered first items had not.

Our work builds on suggestive evidence from a prior study that predictive value may influence subsequent memory (Sherman and Turk-Browne, 2020). This prior study included behavioral and fMRI experiments, whereas the current study used iEEG. Neural measures are an important advance over behavior alone because they can assay predictive representations during passive viewing at encoding. iEEG is superior to fMRI for this purpose because neural activity is sampled at much greater temporal resolution and activity reflects instantaneous electrical potentials rather than hemodynamic responses smoothed and delayed in time. This provides much greater confidence that the upcoming category was being represented before its appearance and thus was truly predictive. Moreover, the prior study showed a negative relationship between prediction and memory across participants, whereas the current study established this relationship within participant. This is also an important advance because an across-participant relationship does not provide strong evidence for the claim that prediction during encoding impairs memory. Such a relationship could reflect generic individual differences such that, for example, a participant with better overall memory generates the same weak prediction on both remembered and forgotten trials. In contrast, in this study, we were able to link prediction to successful versus unsuccessful memory formation across items. This more sensitive approach yielded other findings not observed in the prior study, including that memory for B items had an opposite, positive relationship with prediction of B. Together, these results provide mechanistic insight into the interaction between predictive value and memory, and speak to theoretical questions about the representations underlying statistical learning and episodic memory.

Nature of representational changes

Several fMRI studies have shown that statistical and related forms of learning can change neural representations of associated items throughout the human brain (Schapiro et al., 2012, 2013; Schlichting et al., 2015; Deuker et al., 2016; Tompary and Davachi, 2017). For example, if exposed to sequential pairs embedded in a continuous stream of objects (akin to the category pairs in the current study), the two objects in a pair come to elicit more similar patterns of fMRI activity from before to after learning, when presented on their own, in the medial temporal lobe cortex and hippocampus (Schapiro et al., 2012). Such integration could be interpreted as evidence that the representations of the paired items merged into a single “unitized” representation of the pair that can be evoked by either item (Fujimichi et al., 2010). Alternatively, the paired items may remain distinct but become associated, such that either can be reactivated by the other through spreading activation (Schapiro et al., 2017). A key difference between these accounts is the timing of how learned representations emerge when one of the items is presented: the merging account predicts that the (same) unitized representation is evoked immediately by

either paired item, whereas the associative account predicts that the presented item is represented immediately while the paired item is represented gradually over time through reactivation. These dynamics cannot be distinguished by fMRI because of its slow temporal resolution, but our iEEG approach may shed light.

On the surface, the results of our frequency tagging analysis may seem to suggest a merged representation of the category pairs. The reliable peak in coherence at the frequency of two consecutive stimuli may suggest that electrodes in visual cortex represented the paired categories as a single unit (Batterink and Paller, 2017). However, the results of our pattern similarity analysis are more consistent with an association between the paired categories. Although we found that both categories in a pair could be represented at the same time (i.e., predictive B evidence during the A Pre trial and lingering A evidence during the B Post trial, relative to no such evidence on X trials), these representations were offset in time. The representation of the A category was robust during both the ON and ISI epochs of the A trial, whereas the representation of the B category was not reliable during the ON epoch and only emerged during the ISI epoch. Thus, our results are more consistent with an associative account in visual cortex. It remains possible that the hippocampus or other brain structures represent statistical regularities through unitized representations. Moreover, one limitation of our study is that we did not measure representations of individual categories before and after learning to directly assess representational change. Although we could not directly measure representational change from before to after learning, we did correlate the category templates measured after learning. Unitization of paired categories would be reflected in increased pattern similarity among paired, relative to unpaired and random categories. We did not find reliable evidence of such representational merging, inconsistent with a unitization account. However, prior studies focused on the unitization of paired items rather than categories. Thus, if we had found evidence of representational merging of paired categories in the current study, it would be unclear whether this reflects unitization in the same way or a qualitatively different kind of representational change.

Predictive interference on memory encoding

The time course of predictive representations also sheds light on the temporal dynamics of the interaction between episodic memory and statistical learning. When examining the overall effect of prediction, we found reliable B evidence during the ISI epoch of A, immediately preceding the appearance of B. However, this result was obtained by averaging across all trials, both remembered and forgotten. Thus, it was possible that, when separated out by subsequent memory, a different pattern would emerge. One possibility is that B evidence would come online earlier for forgotten items, which might suggest that the observed impairment in A memory resulted from interference with perceptual processing of A. To the contrary, the difference in B evidence for remembered versus forgotten A items was clearest during the ISI after A was removed from the screen, which suggests that prediction may interfere with later, post-perceptual stages of processing to impair encoding.

Interestingly, evidence for the current A category was comparable across remembered and forgotten A items. Thus, in this paradigm, variance in memory was explained solely by prediction of the upcoming category, not the strength of perceptual

processing of the category being encoded (Kuhl et al., 2012) nor modulation of this processing by prediction (both of which would have affected A evidence). The lack of a relationship between A evidence and A memory may reflect a trade-off: category evidence may reflect representation of the most diagnostic features of a category, which would enhance memory for these features while impairing memory for idiosyncratic features of particular exemplars. A related account may explain why predictive B evidence was positively linked to B memory (Smith et al., 2013; Thavabalasingam et al., 2016): B evidence during the A ISI may potentiate the diagnostic features of the B category, enhancing the salience of idiosyncratic features of B when it appears to strengthen episodic memory for B. Future studies could test these possibilities by using a more continuous measure of memory precision and by testing on modified items that retain category-diagnostic versus idiosyncratic features.

Our finding that prediction relates to better memory for predictable B items contrasts with findings of enhanced encoding for unpredictable/unexpected items (G. Kim et al., 2014; Greve et al., 2017; Bein et al., 2021). These seemingly divergent findings are difficult to reconcile because predictions in our study were never violated: in the Structured condition, the A in each pair was followed deterministically by B; in the Random condition, although each X was unexpected to some degree, they did not violate a learned expectation. Thus, it is possible that replacing the expected B with another category would have led to even better memory encoding. That said, one interpretation of our finding of enhanced (predictable) B memory that would be consistent with a benefit of prediction error for episodic memory could be that features idiosyncratic to a particular B exemplar (needed to later retrieve this specific episodic memory) may have violated a category-level expectation grounded in the diagnostic (i.e., nonidiosyncratic) features of a category shared across its exemplars. This question, as well as questions above about how the category-level nature of the prediction may have affected memory for A, could be informed by future studies examining effects of item-level prediction on memory.

This work builds on existing theories considering the complex interplay between memory encoding and memory retrieval. To the extent that prediction from statistical learning can be considered associative retrieval (Hindy et al., 2016; Kok and Turk-Browne, 2018), our findings converge with the notion that the brain cycles between mutually exclusive encoding and retrieval states (Hasselmo et al., 2002; Duncan et al., 2012; Long and Kuhl, 2019; Bein et al., 2020), organized by the hippocampal theta cycle (Kerrén et al., 2018; Pacheco Estefan et al., 2021). Further, a recent computational model suggests that predictive uncertainty determines when memories should be encoded or retrieved (Lu et al., 2022). The model accounts for findings that familiar experiences are more likely to evoke retrieval (Patil and Duncan, 2018), and thus may help to explain why predictions from statistical learning are prioritized over episodic encoding.

Neural source of predictions

The current study sought to decode evidence of visual categories and so focused on electrode contacts in visual cortex. This adds to a growing literature on predictive signals in visual cortex (De Lange et al., 2018; H. Kim et al., 2020; Clarke et al., 2022). Importantly, in our previous fMRI study (Sherman and Turk-Browne, 2020), we found evidence of prediction only in the hippocampus. We interpreted the lack of an effect in visual cortex in light of the fact that we were measuring prediction (of B) while other items (A) were being perceived; thus, if visual cortex

preferentially represents onscreen, perceived information, we may not have been sensitive to a weaker, simultaneous prediction effect. Indeed, other fMRI studies have found predictions in visual cortex during the absence or omission of perceptual input (Hindy et al., 2016; Clarke et al., 2022). Using a time-resolved measure like iEEG in the current study provided another solution to this problem, by allowing us to isolate short ON versus ISI time periods when there was versus was not a competing stimulus present, respectively (which fMRI would have been unable to separate). Indeed, we found evidence for prediction during the ISI after the predictive item but not while the predictive item was on the screen. This increased sensitivity to prediction specifically during the ISI period may have also provided a clean enough prediction signal to detect a trial-level relationship with memory.

Although we observe these predictive signals in visual cortex, these signals may originate elsewhere in the brain. A strong candidate is the hippocampus and surrounding medial temporal lobe cortex. In addition to representing predictions (Kok and Turk-Browne, 2018; Sherman and Turk-Browne, 2020; Reddy et al., 2021), the hippocampus interfaces between perception and memory (Treder et al., 2021) and has been shown to drive reinstatement of predicted information in visual cortex (Bosch et al., 2014; Tanaka et al., 2014; Hindy et al., 2016; Danker et al., 2017).

Beyond generating predictions, the hippocampus may also be the nexus of the interaction between episodic memory and statistical learning, given its fundamental role in both functions (Schapiro et al., 2017). Indeed, given the necessity of the hippocampus for episodic memory, our study raises questions about how the representations of perceived and predicted categories in visual cortex are routed into the hippocampus for encoding. One intriguing possibility is that these representations are prioritized according to biased competition (Desimone, 1998; Hutchinson et al., 2016), leading to preferential routing and subsequent encoding of predicted, but not perceived, information in the hippocampus. Relatedly, recent work had found that encoding versus retrieval states are associated with distinct patterns of activity in visual cortex (Long and Kuhl, 2021), suggesting that representations in visual regions may be fundamentally shaped by memory state in the hippocampus.

The patients in the current study had relatively few contacts in the hippocampus and medial temporal lobe cortex, precluding careful analysis of prediction in these regions and how it relates to visual cortex. Future studies with a larger cohort of patients and/or high-density hippocampal recordings would be useful for this purpose. Such studies could also provide a more direct link between statistical learning-based prediction and encoding/retrieval modes by examining how hippocampal theta phase (Kerrén et al., 2018; Pacheco Estefan et al., 2021) relates to predictive signals in visual cortex. Likewise, future studies could disrupt the hippocampus through stimulation to establish its causal role in predictive representations in visual cortex.

Limitations of the current study

In the current study, we exploited the high signal-to-noise of intracranial recordings in a small sample of patients. Motivated by the ability to densely sample neural data within this rare population, we focused our experimental design on optimizing neural measures. This led to a few limitations.

Our primary evidence of statistical learning came from neural rather than behavioral measures, namely, neural entrainment at

the pair frequency and category prediction in pattern similarity. We did not have any direct behavioral measures of statistical learning, such as faster response times for predictable items during learning (Gómez et al., 2011; Siegelman et al., 2018) or familiarity judgments about regularities after learning (Fiser and Aslin, 2002; Turk-Browne et al., 2005; Brady and Oliva, 2008). We could not assess statistical learning behaviorally during the encoding phase because we used passive viewing (to reduce task complexity for patients) and because the images were presented too rapidly for manual responses (to enable neural measures of entrainment). We did not include a separate behavioral test of statistical learning after the encoding phase because of limited testing time with the patients that required us to prioritize the neural measures and the behavioral memory test most central to the hypothesis. Future work should consider relating neural signatures of statistical learning from iEEG to more direct behavioral measures of statistical learning, as has been done with scalp EEG (Batterink and Paller, 2017) and fMRI (Karuza et al., 2013).

Statistical learning was also measured indirectly via performance on the recognition memory test. We found reduced memory for predictive A items in the episodic memory test, a replication of prior work (Sherman and Turk-Browne, 2020). This effect provides some evidence of learning because the pairs were novel and arbitrary; and thus, A was only predictive (of B) as a result of new learning. Given that the only difference between A and X was the added predictiveness of A; reduced memory for A relative to X therefore must reflect this learning. That said, there are some limitations to this behavioral effect. Specifically, it was present only in hit rate for A (saying “old” to old exemplars), and not in A', a measure of sensitivity that corrects for false alarm rate for A (saying “old” to new exemplars). The lack of an A' effect resulted from a trend toward lower false alarm rates for A than X. Such a result could suggest a criterion shift for A items (less likely to say “old” in general). However, the prior study (Sherman and Turk-Browne, 2020), which had more statistical power, did not find a similar trend in false alarm rates; rather, there was a similar trend across hit rate and A'. Furthermore, the fact that Structured and Random conditions were presented in separate blocks in the current study (to enable frequency tagging) as opposed to intermixed in the prior study complicates the interpretation of weaker differences between A and X, as they could be confounded with time-dependent differences in the patients' motivation, attention, and/or symptoms. Nevertheless, we were able to leverage variance in memory within A items of the Structured condition, by relating memory to trial-by-trial neural prediction.

Last, we adopted a subblock structure, in which individual exemplars repeated 4 times before switching to new exemplars (but holding the category pairs constant). This choice was made to balance the rapid presentation of stimuli needed for the neural frequency tagging analyses with providing sufficient exposure to the images so that some would be later remembered. Although we found some evidence that neural entrainment to the pairs increased across Structured subblocks, there was little evidence of a learning trajectory in the behavioral or predictive neural measures. It is possible that exemplar repetition in the subblocks may have allowed learning to asymptote after only one or a few subblocks (Turk-Browne et al., 2009), eliminating the possibility of finding a more gradual change in these measures across subblocks. These analyses are further limited by the small number of patients relative to prior work with healthy individuals that found clearer learning effects in behavior (Sherman and Turk-

Browne, 2020). Future studies could tailor their experimental designs to optimize detection of a learning trajectory, for example, by foregoing neural entrainment and presenting images once for a longer duration or by introducing more complex regularities.

In conclusion, in examining the trade-off between prediction and memory encoding, our work suggests a novel theoretical perspective on why predictive value shapes memory. We argue that, because memory is capacity- and resource-limited, memory systems must prioritize which information to encode. When prior statistical learning enables useful prediction of an upcoming experience, that prediction takes precedence over encoding. In this way, encoding is focused adaptively on experiences for which there is room to develop stronger predictions.

References

- Aitken F, Kok P (2022) Hippocampal representations switch from errors to predictions during acquisition of predictive associations. *Nat Commun* 13:1–13.
- Aitken F, Turner G, Kok P (2020) Prior expectations of motion direction modulate early sensory processing. *J Neurosci* 40:6389–6397.
- Aly M, Turk-Browne NB (2017) How hippocampal memory shapes, and is shaped by, attention. In: *The hippocampus from cells to systems*, pp 369–403. New York: Springer.
- Batterink LJ, Paller KA (2017) Online neural monitoring of statistical learning. *Cortex* 90:31–45.
- Bein O, Duncan K, Davachi L (2020) Mnemonic prediction errors bias hippocampal states. *Nat Commun* 11:1–11.
- Bein O, Plotkin NA, Davachi L (2021) Mnemonic prediction errors promote detailed memories. *Learn Mem* 28:422–434.
- Biderman N, Bakkour A, Shohamy D (2020) What are memories for? the hippocampus bridges past experience with future decisions. *Trends Cogn Sci* 24:542–556.
- Bosch SE, Jeehe JF, Fernández G, Doeller CF (2014) Reinstatement of associative memories in early visual cortex is signaled by the hippocampus. *J Neurosci* 34:7493–7500.
- Brady TF, Oliva A (2008) Statistical learning using real-world scenes: extracting categorical regularities without conscious intent. *Psychol Sci* 19:678–685.
- Choi D, Batterink LJ, Black AK, Paller KA, Werker JF (2020) Preverbal infants discover statistical word patterns at similar rates as adults: evidence from neural entrainment. *Psychol Sci* 31:1161–1173.
- Clarke A, Crivelli-Decker J, Ranganath C (2022) Contextual expectations shape cortical reinstatement of sensory representations. *J Neurosci* 42:5956–5965.
- Cowan ET, Schapiro AC, Dunsmoor JE, Murty VP (2021) Memory consolidation as an adaptive process. *Psychon Bull Rev* 28:1796–1810.
- Danker JF, Tompary A, Davachi L (2017) Trial-by-trial hippocampal encoding activation predicts the fidelity of cortical reinstatement during subsequent retrieval. *Cereb Cortex* 27:3515–3524.
- De Brigard F (2014) Is memory for remembering? Recollection as a form of episodic hypothetical thinking. *Synthese* 191:155–185.
- De Lange FP, Heilbron M, Kok P (2018) How do expectations shape perception? *Trends Cogn Sci* 22:764–779.
- Demarchi G, Sanchez G, Weisz N (2019) Automatic and feature-specific prediction-related neural activity in the human auditory system. *Nat Commun* 10:1–11.
- Desimone R (1998) Visual attention mediated by biased competition in extrastriate visual cortex. *Philos Trans R Soc Lond B Biol Sci* 353:1245–1255.
- Deuker L, Bellmund JL, Schröder TN, Doeller CF (2016) An event map of memory space in the hippocampus. *Elife* 5:e16534.
- Dickerson KC, Adcock RA (2018) Motivation and memory. In: *Stevens' Handbook of experimental psychology and cognitive neuroscience*. John Wiley & Sons, Inc. <https://doi.org/10.1002/9781119170174.epcn107>.
- Dolcos F, Katsumi Y, Weymar M, Moore M, Tsukiura T, Dolcos S (2017) Emerging directions in emotional episodic memory. *Front Psychol* 8:1867.

- Duncan K, Sadanand A, Davachi L (2012) Memory's penumbra: episodic memory decisions induce lingering mnemonic biases. *Science* 337:485–487.
- Efron B, Tibshirani R (1986) Bootstrap methods for standard errors, confidence intervals, and other measures of statistical accuracy. *Statist Sci* 1:54–75.
- Ekman M, Kok P, de Lange FP (2017) Time-compressed preplay of anticipated events in human primary visual cortex. *Nat Commun* 8:1–9.
- Endress AD, Johnson SP (2021) When forgetting fosters learning: a neural network model for statistical learning. *Cognition* 213:104621.
- Fiser J, Aslin RN (2002) Statistical learning of higher-order temporal structure from visual shape sequences. *J Exp Psychol Learn Mem Cogn* 28:458–467.
- Fujimichi R, Naya Y, Koyano KW, Takeda M, Takeuchi D, Miyashita Y (2010) Unitized representation of paired objects in area 35 of the macaque perirhinal cortex. *Eur J Neurosci* 32:659–667.
- Gebhart AL, Aslin RN, Newport EL (2009) Changing structures in mid-stream: learning along the statistical garden path. *Cogn Sci* 33:1087–1116.
- Goldfarb EV (2019) Enhancing memory with stress: progress, challenges, and opportunities. *Brain Cogn* 133:94–105.
- Gómez DM, Bion RA, Mehler J (2011) The word segmentation process as revealed by click detection. *Lang Cogn Processes* 26:212–223.
- Greve A, Cooper E, Kaula A, Anderson MC, Henson R (2017) Does prediction error drive one-shot declarative learning? *J Mem Lang* 94:149–165.
- Grier JB (1971) Nonparametric indexes for sensitivity and bias: computing formulas. *Psychol Bull* 75:424–429.
- Hasselmo ME, Bodelón C, Wyble BP (2002) A proposed function for hippocampal theta rhythm: separate phases of encoding and retrieval enhance reversal of prior learning. *Neural Comput* 14:793–817.
- Henin S, Turk-Browne NB, Friedman D, Liu A, Dugan P, Flinker A, Doyle W, Devinsky O, Melloni L (2021) Learning hierarchical sequence representations across human cortex and hippocampus. *Sci Adv* 7:eabc4530.
- Hindy NC, Ng FY, Turk-Browne NB (2016) Linking pattern completion in the hippocampus to predictive coding in visual cortex. *Nat Neurosci* 19:665–667.
- Hutchinson JB, Pak SS, Turk-Browne NB (2016) Biased competition during long-term memory formation. *J Cogn Neurosci* 28:187–197.
- Jenkinson M, Bannister P, Brady M, Smith S (2002) Improved optimization for the robust and accurate linear registration and motion correction of brain images. *Neuroimage* 17:825–841.
- Jenkinson M, Beckmann CF, Behrens TE, Woolrich MW, Smith SM (2012) FSL. *Neuroimage* 62:782–790.
- Jenkinson M, Smith S (2001) A global optimisation method for robust affine registration of brain images. *Med Image Anal* 5:143–156.
- Jungé JA, Scholl BJ, Chun MM (2007) How is spatial context learning integrated over signal versus noise? A primacy effect in contextual cueing. *Vis Cogn* 15:1–11.
- Karuza EA, Newport EL, Aslin RN, Starling SJ, Tivarus ME, Bavelier D (2013) The neural correlates of statistical learning in a word segmentation task: an fMRI study. *Brain Lang* 127:46–54.
- Kerrén C, Linde-Domingo J, Hanslmayr S, Wimber M (2018) An optimal oscillatory phase for pattern reactivation during memory retrieval. *Curr Biol* 28:3383–3392.e6.
- Kim G, Lewis-Peacock JA, Norman KA, Turk-Browne NB (2014) Pruning of memories by context-based prediction error. *Proc Natl Acad Sci USA* 111:8997–9002.
- Kim H, Schlichting ML, Preston AR, Lewis-Peacock JA (2020) Predictability changes what we remember in familiar temporal contexts. *J Cogn Neurosci* 32:124–140.
- Kok P, Turk-Browne NB (2018) Associative prediction of visual shape in the hippocampus. *J Neurosci* 38:6888–6899.
- Kok P, Failing MF, de Lange FP (2014) Prior expectations evoke stimulus templates in the primary visual cortex. *J Cogn Neurosci* 26:1546–1554.
- Kok P, Mostert P, de Lange FP (2017) Prior expectations induce prestimulus sensory templates. *Proc Natl Acad Sci USA* 114:10473–10478.
- Kuhl BA, Rissman J, Wagner AD (2012) Multi-voxel patterns of visual category representation during episodic encoding are predictive of subsequent memory. *Neuropsychologia* 50:458–469.
- Long NM, Kuhl BA (2019) Decoding the tradeoff between encoding and retrieval to predict memory for overlapping events. *Neuroimage* 201:116001.
- Long NM, Kuhl BA (2021) Cortical representations of visual stimuli shift locations with changes in memory states. *Curr Biol* 31:1119–1126.e5.
- Lu Q, Hasson U, Norman KA (2022) A neural network model of when to retrieve and encode episodic memories. *Elife* 11:e74445.
- Oostenveld R, Fries P, Maris E, Schoffelen JM (2011) Fieldtrip: open source software for advanced analysis of MEG, EEG, and invasive electrophysiological data. *Comput Intell Neurosci* 2011:156869.
- Pacheco Estefan D, Zucca R, Arsiwalla X, Principe A, Zhang H, Rocamora R, Axmacher N, Verschure PF (2021) Volitional learning promotes theta phase coding in the human hippocampus. *Proc Natl Acad Sci USA* 118:e2021238118.
- Papademetris X, Jackowski MP, Rajeevan N, DiStasio M, Okuda H, Constable RT, Staib LH (2006) BioImage suite: an integrated medical image analysis suite: an update. *Insight J* 2006:209.
- Patil A, Duncan K (2018) Lingering cognitive states shape fundamental mnemonic abilities. *Psychol Sci* 29:45–55.
- Reddy L, Self MW, Zoefel B, Poncet M, Possel JK, Peters JC, Baayen JC, Idema S, VanRullen R, Roelfsema PR (2021) Theta-phase dependent neuronal coding during sequence learning in human single neurons. *Nat Commun* 12:1–9.
- Schacter DL, Addis DR, Hassabis D, Martin VC, Spreng RN, Szpunar KK (2012) The future of memory: remembering, imagining, and the brain. *Neuron* 76:677–694.
- Schapiro AC, Kustner LV, Turk-Browne NB (2012) Shaping of object representations in the human medial temporal lobe based on temporal regularities. *Curr Biol* 22:1622–1627.
- Schapiro AC, Rogers TT, Cordova NI, Turk-Browne NB, Botvinick MM (2013) Neural representations of events arise from temporal community structure. *Nat Neurosci* 16:486–492.
- Schapiro AC, Turk-Browne NB, Botvinick MM, Norman KA (2017) Complementary learning systems within the hippocampus: a neural network modelling approach to reconciling episodic memory with statistical learning. *Philos Trans R Soc Lond B Biol Sci* 372:20160049.
- Schlichting ML, Mumford JA, Preston AR (2015) Learning-related representational changes reveal dissociable integration and separation signatures in the hippocampus and prefrontal cortex. *Nat Commun* 6:1–10.
- Sherman BE, Turk-Browne NB (2020) Statistical prediction of the future impairs episodic encoding of the present. *Proc Natl Acad Sci USA* 117:22760–22770.
- Sherman BE, Graves KN, Turk-Browne NB (2020) The prevalence and importance of statistical learning in human cognition and behavior. *Curr Opin Behav Sci* 32:15–20.
- Siegelman N, Bogaerts L, Kronenfeld O, Frost R (2018) Redefining 'learning' in statistical learning: what does an online measure reveal about the assimilation of visual regularities? *Cogn Sci* 42:692–727.
- Smith TA, Hasinski AE, Sederberg PB (2013) The context repetition effect: predicted events are remembered better, even when they don't happen. *J Exp Psychol Gen* 142:1298–1308.
- Tanaka KZ, Pevzner A, Hamidi AB, Nakazawa Y, Graham J, Wiltgen BJ (2014) Cortical representations are reinstated by the hippocampus during memory retrieval. *Neuron* 84:347–354.
- Thavabalasingam S, O'Neil EB, Zeng Z, Lee AC (2016) Recognition memory is improved by a structured temporal framework during encoding. *Front Psychol* 6:2062.
- Tompary A, Davachi L (2017) Consolidation promotes the emergence of representational overlap in the hippocampus and medial prefrontal cortex. *Neuron* 96:228–241.e5.
- Treder MS, et al. (2021) The hippocampus as the switchboard between perception and memory. *Proc Natl Acad Sci USA* 118:e2114171118.
- Turk-Browne NB, Jungé JA, Scholl BJ (2005) The automaticity of visual statistical learning. *J Exp Psychol Gen* 134:552–564.
- Turk-Browne NB, Scholl BJ, Chun MM, Johnson MK (2009) Neural evidence of statistical learning: efficient detection of visual regularities without awareness. *J Cogn Neurosci* 21:1934–1945.
- Walther DB, Caddigan E, Fei-Fei L, Beck DM (2009) Natural scene categories revealed in distributed patterns of activity in the human brain. *J Neurosci* 29:10573–10581.