

Research Articles: Behavioral/Cognitive

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

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

Cite as: J. Neurosci 2022; 10.1523/JNEUROSCI.0708-22.2022

Received: 7 April 2022

Revised: 10 October 2022

Accepted: 13 October 2022

This Early Release article has been peer-reviewed and accepted, but has not been through the composition and copyediting processes. The final version may differ slightly in style or formatting and will contain links to any extended data.

Alerts: Sign up at www.jneurosci.org/alerts to receive customized email alerts when the fully formatted version of this article is published.

1 Temporal dynamics of competition 2 between statistical learning and 3 episodic memory in intracranial 4 recordings of human visual cortex

5 Brynn E. Sherman¹, Kathryn N. Graves¹, David M. Huberdeau¹, Imran H. Quraishi²,
6 Eyiymisi C. Damisah³, Nicholas B. Turk-Browne^{1,4}

7 ¹Department of Psychology, Yale University; ²Department of Neurology, Yale University;

8 ³Department of Neurosurgery, Yale University; ⁴Wu Tsai Institute, Yale University

9

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

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

32

33 **Introduction**

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

39 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
40 likely to happen when and where (**De Brigard, 2014; Biderman et al., 2020**). This knowledge is the result
41 of statistical learning, which identifies patterns or regularities in the environment that repeat over time
42 (**Sherman et al., 2020; Endress and Johnson, 2021**) and form the basis of predictions (**De Lange et al., 2018**).
43 We hypothesize that the availability of these predictions during encoding affects whether a new memory is
44 formed. Namely, if one of the main objectives of long-term memory is to enable prediction, in the service of
45 adaptive behavior, experiences that already generate a prediction may not need to be encoded. In contrast,
46 experiences that yield uncertainty about what will happen next may be more important to store because
47 they can help learn over time what should have been expected. Note that this is distinct from whether
48 an experience being encoded was itself expected or unexpected, which also affects subsequent memory
49 (**Greve et al., 2017; Bein et al., 2021**); rather, we argue that the process of generating a prediction based on
50 the experience impedes its encoding.

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

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

68 Additionally, in our prior work, we only found a relationship between prediction and encoding across
69 participants. Average fMRI evidence for the category of second items during first items was negatively
70 associated with overall memory performance for first items. However, this could reflect a generic individual
71 difference — that participants who make more predictions tend to have worse memory — rather than
72 prediction having a mechanistic effect on encoding. According to the latter account, whether a participant
73 remembers or forgets a given item should depend on whether that item triggered a prediction during its
74 encoding. This requires testing for a relationship between prediction and encoding across items within
75 participant. Time-resolved methods with denser sampling of individual trials could better enable trial-level
76 estimates of prediction necessary for within-participant subsequent memory analyses.

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

90 To track statistical learning in the brain, we employed 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.

99 Although the hippocampus may be the nexus of competition between statistical prediction and episodic
100 encoding (**Schapiro et al., 2017; Sherman and Turk-Browne, 2020**), hippocampal signals may be relayed and
101 reinstated throughout the cortical hierarchy (**Bosch et al., 2014; Tanaka et al., 2014; Danker et al., 2017;**
102 **Hindly et al., 2016; Aitken and Kok, 2022; Clarke et al., 2022**) and frequency tagging (**Henin et al., 2021**) in vi-
103 sual cortex. This allowed us to test our hypotheses robustly in epilepsy patients whose clinical care resulted
104 in extensive electrode coverage in visual cortex but not the hippocampus.

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

108 Materials and Methods

109 Participants

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

115 A summary of patient demographics, clinical details, and research participation can be found in **Table 1**.
116 Given electrode coverage and usable data, we retained 9 patients in the behavioral analyses, 8 patients in
117 the neural frequency tagging analyses, and 7 patients in the neural category evidence analyses.

118 iEEG recordings

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

124 iEEG preprocessing

125 iEEG preprocessing was carried out in FieldTrip (**Oostenveld et al., 2011**). A notch filter was applied to re-
126 move 60-Hz line noise. No re-referencing was applied, except for one patient, whose reference was in visual
127 cortex, resulting in a visual-evoked response in all electrodes; for this patient, we re-referenced the data to a
128 white matter contact in the left anterior cingulate cortex. Data were downsampled to 256 Hz and segmented
129 into trials using the triggers.

130 Electrode selection

131 Patients’ electrode contact locations were identified using their post-operative CT and MRI scans. Recon-
132 structions were completed in BioImage Suite (**Papademetris et al., 2006**) and were subsequently registered
133 to the patient’s pre-operative MRI scan, resulting in contact locations projected into the patient’s pre-operative
134 space. The resulting files were converted from the Bioimagesuite format (.MGRID) into native space coordi-
135 nates using FieldTrip functions. The coordinates were then used to create a region of interest (ROI) in FSL
136 (**Jenkinson et al., 2012**), with the coordinates of each contact occupying one voxel in the mask (**Figure 2**).

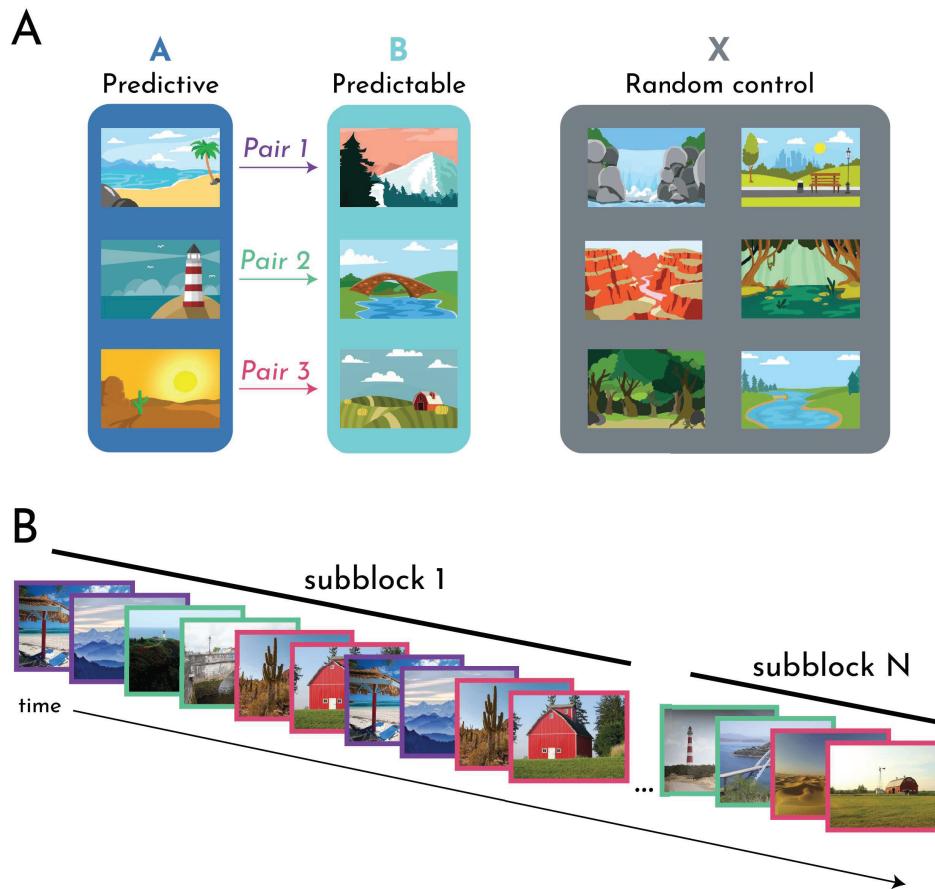


Figure 1. Task design. (A) Example scene category pairings for one participant. Three of 12 categories were assigned to condition A. Each A category was reliably followed by one of three other categories assigned to condition B to create pairs. The remaining six categories assigned to condition X were 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 four times per subblock, and new exemplars were introduced for the next subblock. For the Structured block, the category pairs remained consistent across subblocks. Category pairs are denoted by a colored frame, corresponding to the A-B pairs (and colored arrows) in subpanel A.

Table 1. Patient Information.

ID	Age	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	-

Description of patient participation. ID: patient participation number. Age: in years. Sex: M = Male, F = Female. 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).

137 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 2mm standard brain by combining the Occipital Lobe ROI from the MNI Structural Atlas and the following ROIs from the Harvard-Oxford
 138 Cortical Structural Atlas: Inferior Temporal Gyrus (temporooccipital part), Lateral Occipital Cortex (superior
 139 division), Lateral Occipital Cortex (inferior division), Intracalcarine Cortex, Cuneal Cortex, Parahippocampal
 140 Gyrus (posterior division), Lingual Gyrus, Temporal Occipital Fusiform Cortex, Occipital Fusiform Gyrus,
 141 Supracalcarine Cortex, Occipital Pole. Each ROI was thresholded at 10% and then concatenated together to
 142 create a single mask of visual cortex.
 143

144 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 pre-operative
 145 anatomical scan to the MNI T1 2mm standard brain template using linear registration (FSL FLIRT (*Jenkinson and Smith,*
 146 *2001; Jenkinson et al., 2002*)) with 12 degrees of freedom. This registration was then inverted and used to
 147 bring the visual cortex mask into each participant's native space.
 148

149 In order to ensure that the visual cortex mask captured the anatomical areas we intended, we manually
 150 assessed its overlap between the electrodes and made a few manual adjustments to the electrode definition.
 151 For example, due to noise in the registrations between post-operative and pre-operative space, as well
 152 as from pre-operative space and standard space, some grid and strip contacts appeared slightly outside of
 153 the brain, despite being on the surface of the patient's brain. Thus, contacts such as these were included as
 154 "visual" even if they were slightly outside of the bounds of the mask. Additionally, due to the liberal thresh-
 155 olds designed to capture broad visual regions, some portions of the parahippocampal gyrus area contained
 156 the hippocampus. Contacts within mask boundaries but clearly in the hippocampus were excluded.
 157

158 Experimental Design

159 Participants completed the experiment on a MacBook Pro laptop while seated in their hospital bed. The task
 160 consisted of up to four runs: two runs of the Structured block and two runs of the Random block. We aimed
 161 to collect all four runs from each patient, but required a minimum of one run per condition for subject
 162 inclusion. Given that the order of structured vs. random information can impact learning (*Jungé et al.,*
 163 *2007; Gebhart et al., 2009*), the run order was counterbalanced within and across participants (i.e., some
 164 participants received Structured-Random-Random-Structured and others Random-Structured-Structured-

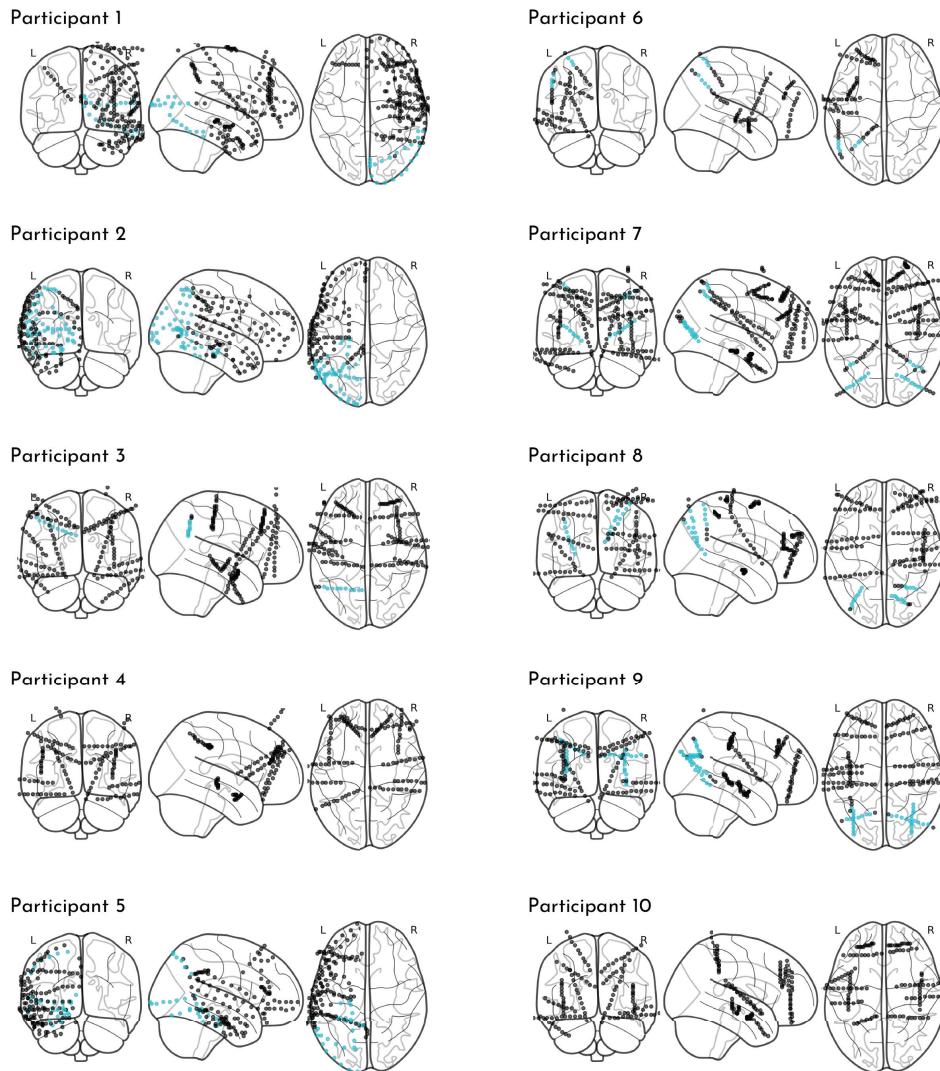


Figure 2. Electrode coverage. The contact locations on the grid, strip, and/or depth electrodes for each participant are plotted as circles in standard brain space. Contacts colored in blue were localized to the visual cortex mask.

¹⁶⁵ Random). Participants completed the runs across 1-3 testing sessions based on the amount of testing time
¹⁶⁶ available between clinical care, family visits, and rest times.

¹⁶⁷ Each run consisted of an encoding phase and a memory phase. During the encoding phase, participants
¹⁶⁸ viewed a rapid stream of scene images, during which they were asked to passively view the scenes. Partici-
¹⁶⁹ pants were told that their memory for the scenes would be tested in order to encourage them to pay close
¹⁷⁰ attention. Each scene was presented for 267 ms, followed by a 267 ms inter-stimulus interval (ISI) period
¹⁷¹ during which a fixation cross appeared in the center of the screen. These short presentation times were cho-
¹⁷² sen to optimize the task for the frequency tagging analyses, which involves measuring neural entrainment
¹⁷³ to stimuli.

¹⁷⁴ Within each run, participants viewed a series of images from a set of six scene categories. There were
¹⁷⁵ six categories in the Structured block, and six other categories in the Random block. In the Structured block,
¹⁷⁶ the scenes categories were paired, such that images from one scene category (A) were always followed by
¹⁷⁷ an image from another scene category (B). Thus, A scenes were *predictive* of the category of the upcoming
¹⁷⁸ B scenes, or stated another way, the category of B scenes was *predictable* given the preceding A scenes.
¹⁷⁹ No scene pairs were allowed to repeat back-to-back in the sequence. In the Random block, all six scene
¹⁸⁰ categories (X) could be preceded or followed by any other scene category, making them neither predictive
¹⁸¹ nor predictable. No individual scene categories were allowed to repeat back-to-back.

¹⁸² In total, participants viewed 16 exemplars from each category within each run. To assist patients with
¹⁸³ remembering these briefly presented images, each individual exemplar was shown four times within a run.
¹⁸⁴ Thus, each run was comprised of 16 "subblocks" during which the same set of six exemplar images was
¹⁸⁵ repeated four times (384 trials total). Within each subblock, the order of the pairs/images was randomized,
¹⁸⁶ with the constraints described above of no back-to-back repetitions. The individual exemplars changed after
¹⁸⁷ each subblock, but the category relations were held constant in the Structured block. Participants were not
¹⁸⁸ informed of these category pairings, and thus had to acquire them through exposure.

¹⁸⁹ At the end of each run, participants completed a memory test. Participants were presented with all
¹⁹⁰ 96 unique images from the encoding phase, intermixed with 24 novel foils from the same categories (4
¹⁹¹ foils/category). Participants first had to indicate whether the image was old, meaning it was just presented
¹⁹² in that run's encoding phase, or new, meaning that they had not seen that image at all during the experiment.
¹⁹³ Following their old/new judgment, participants were asked to indicate their confidence in their response, on
¹⁹⁴ a scale of 1 (very unsure) to 4 (very sure). Participants had up to 6 s to make each old/new and confidence
¹⁹⁵ judgment. We quantified episodic memory performance using A' , a non-parametric measure which takes
¹⁹⁶ into account hit rate (HR) and false alarm rate (FA) (**Grier, 1971**):

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

¹⁹⁷ Frequency tagging analyses

¹⁹⁸ We conducted a phase coherence analysis to identify electrode contacts that entrained to image and pair fre-
¹⁹⁹ quencies (**Henin et al., 2021**). For both Structured and Random blocks, the raw signals were concatenated
²⁰⁰ across runs (if more than one per block type) and then segmented into subblocks comprising 24 trials with
²⁰¹ the four repetitions per exemplar. We then converted the raw signals for each subblock into the frequency
²⁰² domain via fast Fourier transform and computed the phase coherence across subblocks for each electrode
²⁰³ using the formula $R^2 = [\Sigma^N \cos\phi]^2 + [\Sigma^N \sin\phi]^2$. Notably, by computing phase coherence between subblocks,
²⁰⁴ we collapsed over the contribution of individual exemplars that repeated within subblock. In other words,
²⁰⁵ entrainment in this analysis was driven by phase-locking that generalized across exemplars. Phase coher-
²⁰⁶ ence was computed separately for each contact in the visual cortex mask, and we then averaged across
²⁰⁷ contacts within participant. We focused on phase coherence at the frequency of image presentation (534
²⁰⁸ ms/image; 1.87 Hz) and pair presentation (1.07 s/pair; 0.93 Hz).

²⁰⁹ Category evidence analyses

²¹⁰ We employed a multivariate pattern similarity approach to assess the timecourse of category responses.
²¹¹ 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 on-screen category or the upcoming category). The following subsections explain this approach in detail.

Frequency decomposition

We employed a Morlet Wavelet approach to decompose raw signals into time-frequency information (**Figure 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 timecourse 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 timecourses; we then subsequently carved into trials using the triggers, yielding a vector of frequency and contact information at each timepoint within a trial.

Subsequent analyses required that each trial have the same number of timepoints. 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 timepoints 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 (**Figure 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 employed 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 timepoint (each of the 138 timepoints within a trial) and tested each classifier on all 138 timepoints 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 timepoints 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 timepoints very early on in a trial (immediately after image onset) would produce high discrimination between categories. We also reasoned that the best timepoint(s) may differ from participant to participant depending on their electrode coverage. Therefore, we devised a participant-specific timepoint 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 timepoints to assess how well training at every timepoint generalized to all other timepoints 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 timepoint. For each participant, we then computed the

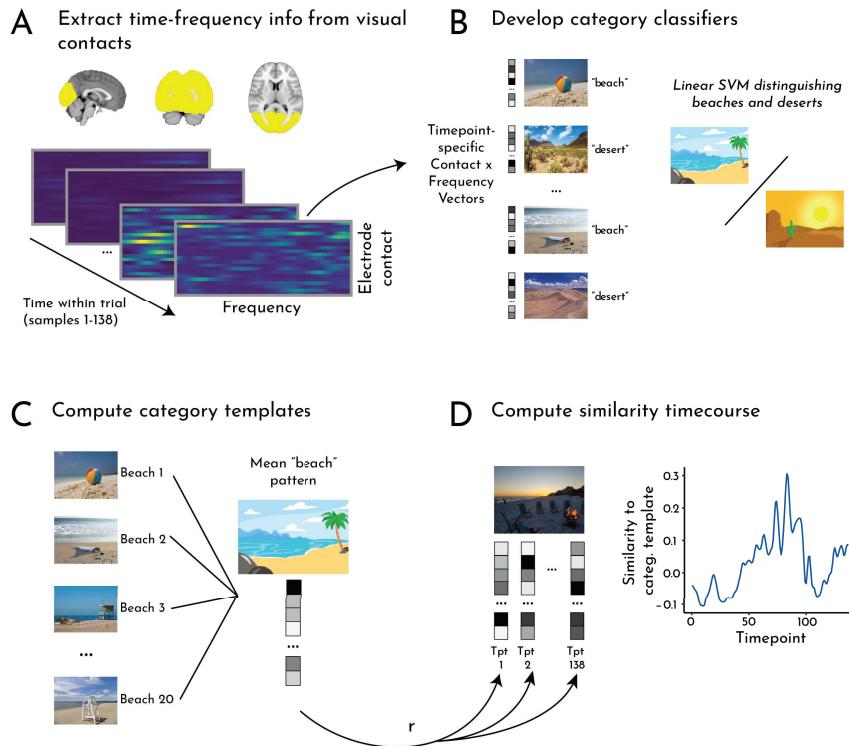


Figure 3. Category evidence analysis pipeline. (A) A Morlet wavelet approach was used to extract time-frequency information from contacts in visual cortex. This resulted in contact by frequency vectors for every timepoint of encoding phase and memory phase trials, which served as the neural patterns 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 timepoint or set of timepoints. The classifiers were then tested on timepoints from held-out data. (C) After a series of feature selection steps, we chose the per-participant top-N timepoint set that produced the best classification accuracy, and then averaged contact by frequency vectors across those timepoints (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 timepoint of each trial/exemplar from that category during the (independent) encoding phase, yielding a timecourse of pattern similarity reflecting neural category evidence.

260 rank order of timepoints with respect to their classification, such that the first ranked timepoint was the
261 one that yielded the highest classification, and the last ranked (138th) timepoint is the one that yielded the
262 lowest classification.

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

270 *Template correlations*

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

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

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

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

295 *Subsequent memory*

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

304 *Alternative classification approaches for feature selection*

305 The category evidence analyses described above rely on a set of binary classifiers trained to distinguish the
306 categories in a given condition (i.e., all combinations of As and Bs in the Structured condition and Xs in the
307 Random condition). However, this approach may lead to interpretational issues. For example, from a binary

308 classifier trained to distinguish two categories (e.g., A1 vs. B1), it is difficult to know whether evidence for
 309 one category (e.g., A1) reflects the presence of that category (A1) or the absence of the other category (B1).
 310 Thus, we replicated all of the above analyses using two alternative approaches.

311 First, we trained a 6-way classifier to distinguish among all six categories of a given condition (A1-A2-
 312 A3-B1-B2-B3 for Structured and X1-X2-X3-X4-X5-X6 for Random). By including more than two classes, this
 313 approach addresses the concern that classification accuracy could be driven by the presence or absence
 314 of a given category. Second, we retained the binary classification approach but trained classifiers to only
 315 discriminate within the A or B categories. That is, instead of 15 classifiers for A/B combinations, there were
 316 6 classifiers (A1-A2, A1-A3, A2-A3, B1-B2, B1-B3, B2-B3). This approach ensures that classification does not
 317 mix evidence for predictive vs. predicted categories.

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

322 Statistical analysis

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

334 Results

335 Memory behavior

336 We first assessed overall performance in the recognition memory test to verify that participants were able
 337 to encode the images into memory. We computed A', a non-parametric measure of sensitivity (**Grier, 1971**),
 338 from test judgments for items from both Structured and Random blocks. All participants had an A' above the
 339 chance level of 0.5 (mean = 0.68; 95% CI = [0.64, 0.70], *p* < 0.001; **Figure 4A**) indicating reliable memory. This
 340 was driven by a higher hit rate (mean = 0.51) than false alarm rate (mean = 0.32; difference 95% CI = [0.14,
 341 0.23], *p* < 0.001). The proportions of items that were subsequently remembered (hit rate) or forgotten (1-hit
 342 rate, or misses) were roughly matched on average, yielding balanced power for within-subject subsequent
 343 memory analyses.

344 We then assessed how statistical learning affected recognition memory. Based on our prior work (**Sherman and Turk-Browne,**
 345 **2020**), we hypothesized that the hit rate for items from the predictive A categories in the Structured blocks
 346 would be lower than the hit rate for items from the control X categories in the Random blocks. Indeed, we
 347 replicated this key behavioral finding (**Figure 4B**), with a significantly lower hit rate for A (mean = 0.48) than
 348 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
 349 from A (difference 95% CI = [-0.10, 0.059], *p* = 0.51) or X (difference 95% CI = [-0.094, 0.053], *p* = 0.66).

350 The false alarm rate for X (mean = 0.36) was numerically higher than A (mean = 0.28; difference 95%
 351 CI = [-0.0023, 0.16], *p* = 0.064); X was significantly higher than B (mean = 0.29; difference 95% CI = [0.0069,
 352 0.13], *p* = 0.028), though A and B did not differ (difference 95% CI = [-0.074, 0.056], *p* = 0.82). Unlike the
 353 higher hit rate for X than A, which was specifically hypothesized based on prior work, the marginally higher
 354 false alarm rate for X than A was not expected or consistent with previous experiments. Nevertheless, this
 355 complicates interpretation of the hit rate difference as impaired memory for A vs. X. One difference from

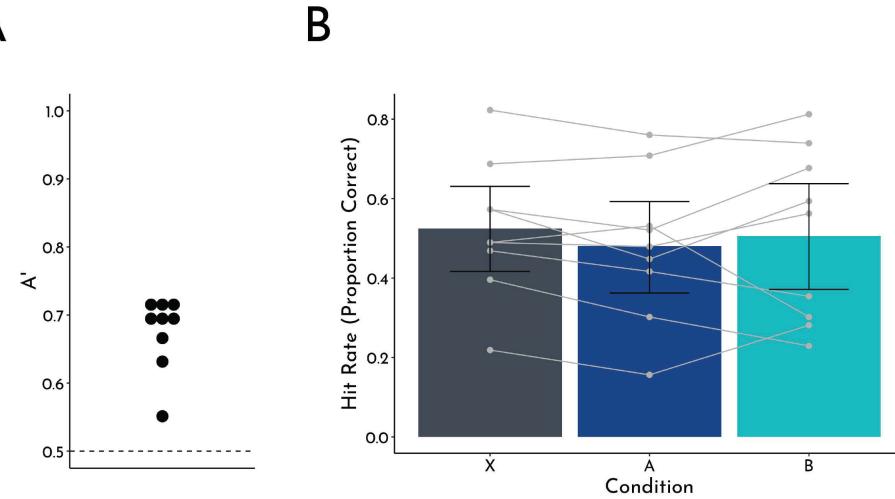


Figure 4. Behavioral results. (A) Overall memory performance collapsed across conditions. A' (a sensitivity measure for recognition memory) is depicted for each participant as a circle. All participants were above chance (0.5). (B) Hit rate as a function of condition (A: predictive; B: predictable; X: control). Group means are plotted as bars, with errors bars representing the bootstrapped 95% confidence interval across participants. Individual participant data are overlaid with the grey circles and lines.

356 the prior study is the blocking of Structured (A,B) and Random (X) categories, which may have allowed for
 357 differences in strategy or motivation between conditions. Nevertheless, the main memory hypotheses in
 358 the current study rest within the A condition (i.e., which A items are remembered vs. forgotten as a function
 359 of B prediction), rather than on overall condition-wide differences with X (or B).

360 We additionally examined the timecourse of these memory effects by sorting the items into subblocks.
 361 If the deficit in memory for A items arises from the predictive value that they confer, we might expect
 362 that this effect will emerge over time as learning occurs (**Sherman and Turk-Browne, 2020**). We focused
 363 this analysis on the first Structured run of the encoding phase for each participant, in order to equate the
 364 amount of data and corresponding opportunity for learning across participants (some had one run, others
 365 two). We quantified change over time for each participant as the Spearman rank correlation of subblock
 366 number with hit rate for A (averaged across items in each subblock), expecting a negative correlation. The
 367 resulting within-participant relationship was not reliable at the group level (mean rho = -0.038; 95% CI =
 368 [-0.27, 0.19], $p = 0.77$). This null effect of a learning trajectory stands in contrast with what we observed in
 369 **Sherman and Turk-Browne (2020)**, perhaps related to the smaller number of participants or differences in
 370 task design (e.g., the use of 'subblocks') in the current study.

371 Neural frequency tagging

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

379 Consistent with our hypotheses and prior work (**Henin et al., 2021**), there were distinct peaks in phase
 380 coherence at both the image and pair frequencies in Structured blocks, but only at the image frequency in
 381 Random blocks (**Figure 5B**). To confirm the reliability of these peaks, we contrasted the coherence at the

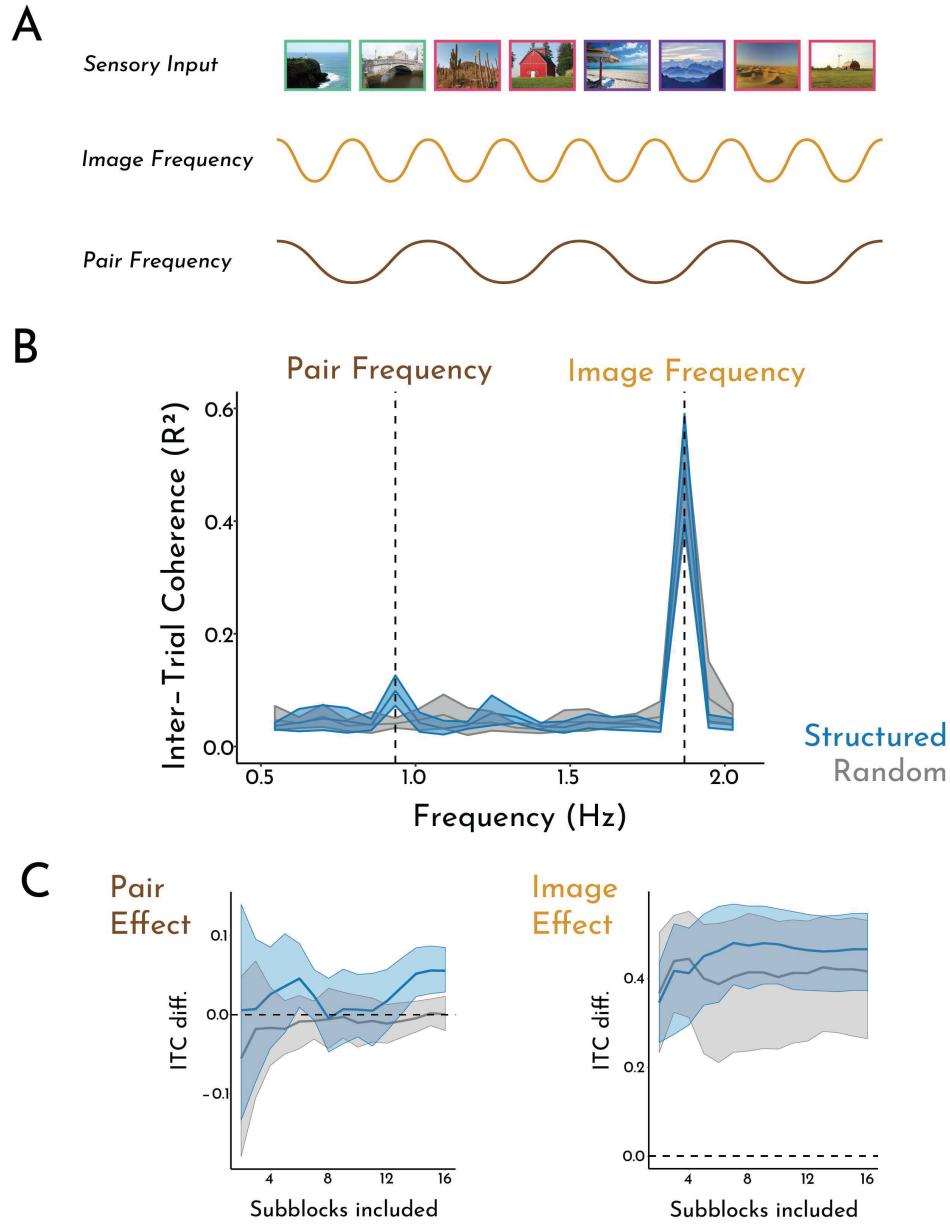


Figure 5. Neural frequency tagging analysis. (A) Schematic of analysis and hypothesized neural oscillations. We expect entrainment of visual contacts at the frequency of images in both blocks. In the Structured block, we also expect entrainment at the frequency of category pairs. (B) These hypotheses 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 confidence intervals across participants.

382 frequency of interest (image: 1.87 Hz; pair: 0.93 Hz) against a baseline of the coherence at frequencies
383 neighboring each of the frequencies of interest (± 0.078 Hz). At the image frequency, there were reliable
384 peaks in both the Structured (mean difference = 0.46; 95% CI = [0.37, 0.55], $p < 0.001$) and Random blocks
385 (mean difference = 0.42; 95% CI = [0.28, 0.52], $p < 0.001$). At the pair frequency, there was a reliable peak
386 in Structured blocks (mean difference = 0.059; 95% CI = [0.035, 0.084]), $p < 0.001$), but not Random blocks
387 (mean difference = -0.0027; 95% CI = [-0.016, 0.0085], $p = 0.68$).

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

394 To measure the emergence of these entrainment effects over time, we computed the coherence over an
395 iteratively increasing number of subblocks (**Henin et al., 2021**). Specifically, we first computed the coherence
396 across the first two subblocks, then the first three, and so on, up to all 16 subblocks. As in the behavioral
397 timecourse analyses, we only included the first 16 subblocks per participant (corresponding to the first run
398 of a given condition) in order to equate the opportunity for learning effects across participants. To quantify
399 neural entrainment, we computed the difference in coherence between the frequency of interest and the
400 two neighboring frequencies (as we did above to establish whether peaks were reliable). We then assessed
401 the reliability of that difference, relative to 0, across participants. We hypothesized that coherence at the pair
402 frequency would emerge over time in the Structured condition, but that coherence at the image frequency
403 would be consistently high, even at early timepoints.

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

409 In contrast to the pair frequency that required learning, the image frequency should be driven by the
410 stimuli and thus present early in both conditions. Indeed, coherence at the image frequency was reli-
411 ably high across all numbers of subblocks, in both the Structured and Random conditions (all $ps < 0.001$;
412 **Figure 5C**, right). This lends credence to the interpretation of increasing coherence at the pair frequency
413 over time as reflecting a trajectory of learning.

414 Given our interpretation that entrainment to the pair frequency reflects statistical learning, and given
415 that we expect our key behavioral effect (impaired memory for predictive A items) to depend on statisti-
416 cal learning, we next asked whether these two effects are related. We calculated this relationship within-
417 participant given the small sample for estimating across-participant relationships. Coherence is necessarily
418 measured across trials, and thus we could not relate entrainment on a given trial to memory for that trial.
419 Instead, we computed coherence across neighboring subblocks and estimated neural entrainment to the
420 pairs as the difference in coherence at the pair frequency from the two adjacent frequencies. We then
421 related this neural measure to average A hit rate within the latter of the two neighboring subblocks, expect-
422 ing a negative relationship (stronger pair entrainment associated with worse A memory). For example, the
423 coherence between subblocks 1 and 2 was used to predict behavioral memory in subblock 2 (memory in
424 subblock 1 was excluded from this analysis). The within-participant relationship between neural entrain-
425 ment to pairs and A memory showed a trend at the group level (mean rho = -0.13; 95% CI = [-0.25, 0.020], p
426 = 0.089), though importantly 6/7 participants showed a negative correlation. We repeated this analysis for
427 the image frequency as a control, and found no relationship between neural entrainment to images and A
428 memory (mean rho = -0.072; 95% CI = [-0.24, 0.087], $p = 0.42$).

429 **Scene category decoding and template creation**

430 The neural frequency tagging for pairs in Structured blocks indicates statistical learning of the pairs. This
431 learning should create predictive value for the items from the A categories, which afford a prediction of

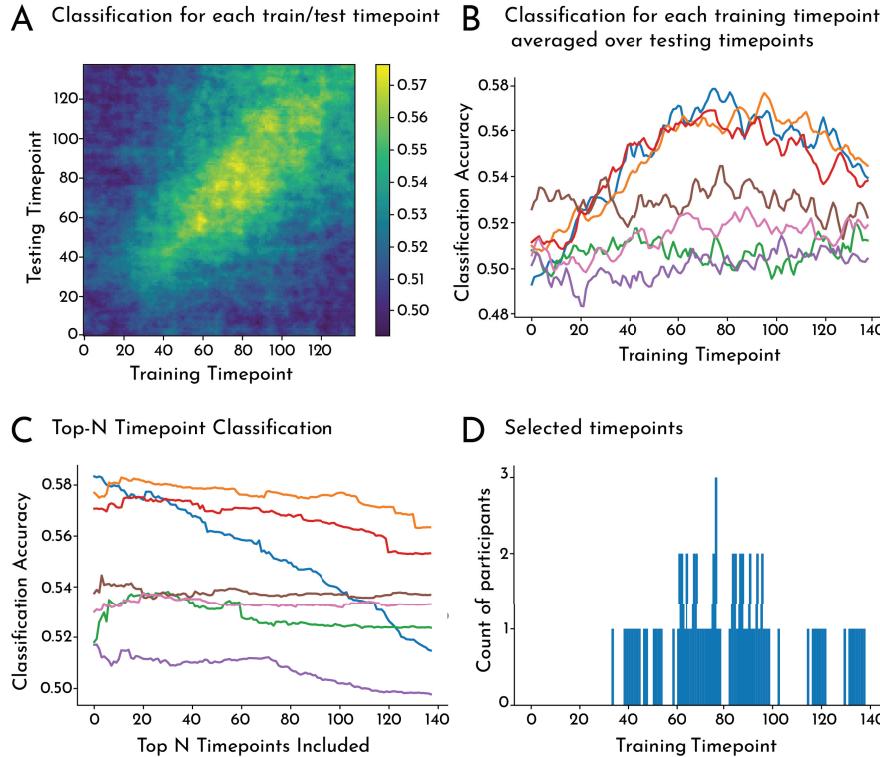


Figure 6. Category decoding and feature selection. (A) To establish overall category decoding accuracy, we trained and tested binary category classifiers separately for all individual timepoints, yielding a temporal generalization matrix. (B) As a first feature-selection step, we computed the average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

432 the associated B category. To test for these predictive representations, we employed a multivariate pattern
 433 similarity approach that extracted neural evidence for visual categories. For each category, we created a
 434 neural template based on the pattern of time-frequency information evoked by each category across visual
 435 contacts. These templates were optimized through a series of steps (described below) for each participant
 436 to ensure maximum category discriminability.

437 First, to verify that the scene categories were indeed discriminable, we developed a series of binary clas-
 438 sifiers to distinguish among the scene categories. Because we were interested in ultimately selecting the
 439 timepoints that produced the best category discrimination, we trained classifiers on a single timepoint (each
 440 of the 138 timepoints within a trial) and tested each classifier on all 138 timepoints at test. **Figure 6A** illus-
 441 trates the classification performance across all of these binary classifiers, averaged across participants. At
 442 the group level (averaging across all train-test combinations), classification performance was above chance
 443 (mean = 0.528; 95% CI = [0.514, 0.542], $p < 0.001$), with each individual participant exhibiting classification
 444 performance greater than the chance level of 0.5.

445 We next aimed to select the training timepoints (per participant) that exhibited the best category dis-
 446 crimination. For each participant and training timepoint, we averaged classification accuracy across all test

447 timepoints (**Figure 6B**). We then ranked the training timepoints by classification accuracy. Next, to find the
 448 set of training timepoints that produced the best classification, we re-ran our classification procedure, but
 449 training on an increasing number of timepoints, starting with the best-performing timepoint, and iteratively
 450 adding timepoints per rank. We then computed the per-participant average classification accuracy for each
 451 set of timepoints **Figure 6C**). Verifying that this feature selection approach worked to optimize category dis-
 452 crimability, we indeed found that using the per-participant top-N timepoints yielded higher classification
 453 accuracy than averaging across all timepoints (mean accuracy = 0.554; 95% CI = [0.536, 0.571], $p < 0.001$);
 454 this was independently true for each participant.

455 We used these per-participant top-N timepoints to create templates of each category. **Figure 6D** illus-
 456 trates the training timepoints which were included in the templates, for one or more participants. To con-
 457 struct the templates, we averaged the contact-by-frequency vectors across the top-N timepoints for all ex-
 458emplars of a given category. We then aimed to quantify the expression of these category templates during
 459 learning (e.g., during the presentation of a predictive A item, is there a representation of the upcoming
 460 B item?). However, given that these templates were created from the memory phase, after learning had
 461 already occurred, it is important to ensure that the templates of paired categories themselves were not
 462 correlated with each other; if so, any effects of prediction during learning could be confounded. At the
 463 group level, the templates of paired categories (e.g., A1-B1) were no more correlated than the templates of
 464 unpaired Structured categories (e.g., A1-B2; mean difference = 0.024; 95% CI = [-0.019, 0.069], $p = 0.30$) or
 465 Random categories (e.g., X1-X2; mean difference = 0.047; 95% CI = [-0.032, 0.127], $p = 0.25$).

466 Category evidence during learning

467 To test for evidence of predictive value, we quantified the expression of these templates in the Structured
 468 and Random blocks. As a check, we expected clear neural evidence for the category of the item being pre-
 469 sented on the screen. Critically, we hypothesized that neural evidence for the upcoming B category would
 470 manifest before its appearance, in response to an A exemplar. We measured these temporal dynamics of
 471 neural category evidence by creating a window of three trials centered on the current item: the trial preced-
 472 ing a trial in which the item appeared ("Pre"), the trial during which the item was on the screen ("Current"),
 473 and the trial succeeding the trial in which the item appeared ("Post"). For example, if category Pair 1 involved
 474 beaches (A1) being followed by mountains (B1), neural evidence for the mountain category was calculated
 475 in response to beach exemplars (Pre), mountain exemplars (Current), and exemplars from the categories
 476 that could appear next in the Structured sequence (A2 or A3 categories). These evidence values were aver-
 477 aged across the categories from the same condition (e.g., B1, B2, and B3 for condition B) and plotted over
 478 time (**Figure 7A**). For statistical analysis, we averaged the neural category evidence for each category across
 479 the timepoints within 6 epochs: when Pre, Current, and Post images were on the screen ("ON") and dur-
 480 ing the fixation period between these trials ("ISI"; **Figure 7B**). We anticipated the evoked response to each
 481 image would span ON and ISI periods (as neural processing of the image would take longer than 267 ms),
 482 but subdividing in this way allowed us to test for the emergence of predictive evidence of B during the ISI
 483 immediately prior to its onset.

484 For Current trials (i.e., the trial when the target category was on screen), we found robust (perceptual)
 485 evidence for both A and B across both the ON epoch (A: mean = 0.0088; 95% CI = [0.0046, 0.013], $p < 0.001$;
 486 B: mean = 0.012; 95% CI = [0.0066, 0.018], $p < 0.001$) and ISI epoch (A: mean = 0.012; 95% CI = [0.0084,
 487 0.015], $p < 0.001$; B: mean = 0.014; 95% CI = [0.0083, 0.019], $p < 0.001$). Neural evidence for X categories
 488 from Random blocks was not reliable during the ON epoch (mean = 0.0046, 95% CI = [-0.00075, 0.012], p
 489 = 0.13) but became robust later in the trial during the ISI epoch (mean = 0.0074; 95% CI = [0.0030, 0.013],
 490 $p < 0.001$). There was greater evidence for B than X categories during both ON (mean difference = 0.0077;
 491 95% CI = [0.00058, 0.015], $p = 0.031$) and ISI epochs (mean difference = 0.0065; 95% CI = [0.00061, 0.012], p
 492 = 0.031). Considering X as a baseline, this difference shows enhanced perceptual processing of predictable
 493 categories. Neural evidence did not differ between A and B categories ($ps > 0.38$) or A and X categories (ps
 494 > 0.28).

495 For Pre trials (i.e., the trial before the target category appeared), we found the hypothesized predictive
 496 neural evidence for the B categories during the ISI epoch (just after its paired A category appeared; mean =

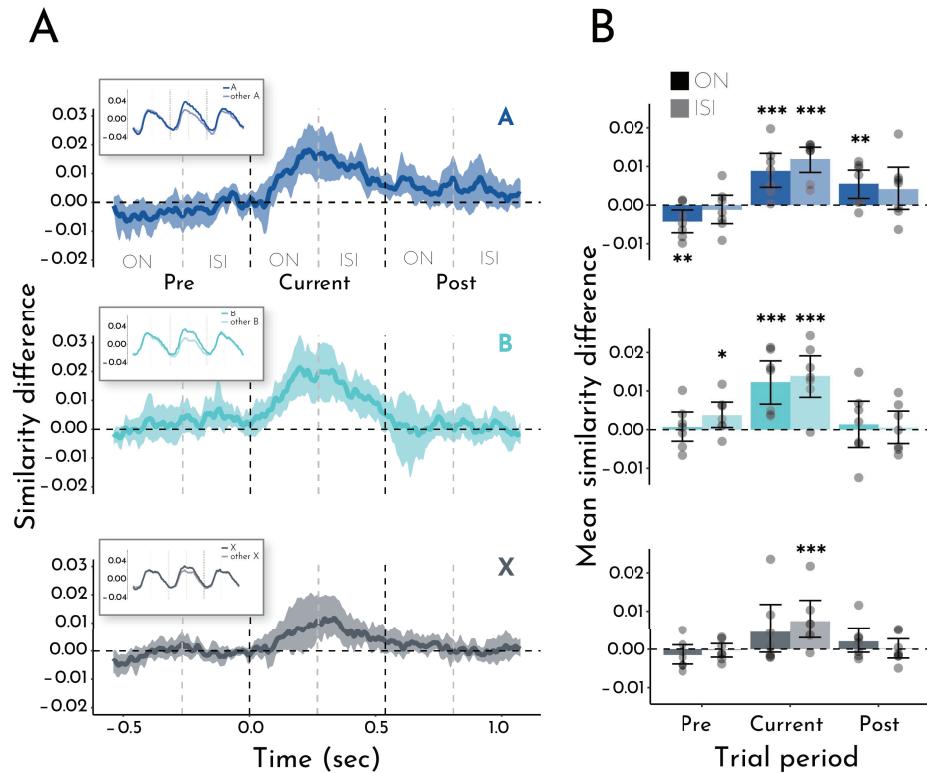


Figure 7. Neural category evidence. (A) Time course of similarity between patterns of neural activity in visual contacts evoked by exemplars from A (predictive), B (predictable), and X (control) categories and category template patterns for A, B, and X, respectively, baselined to average evidence for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

497 0.0037; 95% CI = [0.00054, 0.0071], $p = 0.019$). B evidence was not present during the ON epoch earlier in the
 498 Pre trials (while its paired A category was on screen; mean = 0.00063; 95% CI = [-0.0030, 0.0046], $p = 0.78$);
 499 this may reflect the time needed for associative reactivation of the B category after perceptual processing of
 500 the A item, or anticipation of the timing when B will appear (at the end of the Pre trial). Further supporting
 501 our interpretation that Pre evidence of the B categories reflects prediction, no such evidence was observed
 502 for X during ON (mean = -0.0015; 95% CI = [-0.0039, 0.0012], $p = 0.26$) or ISI epochs (mean = -0.00031; 95%
 503 CI = [-0.0021, 0.0015], $p = 0.73$) or for A during the ISI epoch (mean = -0.0012; 95% CI = [-0.0048, 0.0025], $p =$
 504 0.53). There was *negative* evidence for the upcoming A category during the ON epoch of the Pre trial (mean
 505 = -0.0043; 95% CI = [-0.0072, -0.0013], $p = 0.0052$), but this may have been artifactual (see below). When
 506 contrasting prediction-related signals across conditions, Pre neural evidence for the B categories during
 507 the ISI epoch was reliably greater than X categories (mean difference = 0.0040; 95% CI = [0.00016, 0.0075],
 508 $p = 0.042$) and marginally greater than A categories (mean difference = 0.0049; 95% CI = [-0.00051, 0.010],
 509 $p = 0.075$).

510 For Post trials (i.e., the trial after the target category appeared), we found reliable neural evidence for
 511 the A categories during the ON epoch (i.e., while its paired B category was on screen; mean = 0.0055; 95%
 512 CI = [0.0017, 0.0091], $p = 0.0018$); this effect was not significant during the ISI epoch (mean = 0.0041; 95%
 513 CI = [-0.0011, 0.0098], $p = 0.13$). We did not find Post evidence of B or X categories during either ON or ISI
 514 epochs ($p > 0.80$), nor was Post evidence for A reliably stronger than B or X ($p > 0.16$). Positive evidence of A
 515 during the Post trial may be related to the negative evidence of A during the Pre trial noted above. Because
 516 no back-to-back pair repetitions were allowed, in an A1-B1-A2-B2 trial sequence, A1 and A2 were different
 517 categories. A1 evidence during B1 was considered a Post trial for the A condition, whereas A2 evidence
 518 during B1 was considered a Pre trial for the A condition. Because A1 was one of two baseline categories for
 519 A2 (along with the third A category, A3), Post evidence for A1 during B1 would have been subtracted from
 520 Pre evidence for A2, leading to a negative effect. We tested this by comparing evidence for A2 (Pre) and A1
 521 (Post) during B1 to the neutral A3 only. This weakened the negative Pre evidence for A, during ON (mean =
 522 -0.0027; 95% CI = [-0.0054, 0.00], $p = 0.058$) and ISI epochs (mean = 0.00048; 95% CI = [-0.0022, 0.0038], $p =$
 523 0.82). However, the positive Post evidence for A during the ON epoch remained significant (mean = 0.0081;
 524 95% CI = [0.0036, 0.014], $p < 0.001$).

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

528 First, we created category templates from a 6-way classifier that simultaneously learned to distinguish
 529 the patterns from all categories of a condition. As a check, we first confirmed that this method produced
 530 the same results for Current items. Indeed, as above, we found reliable evidence for both A and B items,
 531 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],
 532 $p < 0.001$) and ISI periods (A: mean = 0.010; 95% CI = [0.0060, 0.014], $p < 0.001$; B: mean = 0.014; 95% CI =
 533 [0.0085, 0.019], $p < 0.001$); evidence for X was reliable during the ISI (mean = 0.0059; 95% CI = [0.0026, 0.0099],
 534 $p < 0.001$), but not ON periods (mean = 0.0037; 95% CI = [-0.0026, 0.012], $p = 0.32$). Critically, we replicated our
 535 key finding of predictive B evidence during the Pre-ISI period (i.e., just after its paired A category appeared;
 536 mean = 0.0035; 95% CI = [0.00042, 0.0066], $p = 0.025$), as well as of lingering A evidence during the Post-ON
 537 period (i.e., while its paired B category was on screen; mean = 0.0049; 95% CI = [0.000059, 0.0095], $p =$
 538 0.049).

539 Second, we retained the binary classification approach but limited the classifiers to category compar-
 540 isons within A or within B, such that the classifiers did not learn to discriminate A vs. B. Although we ex-
 541 pected that this approach would reduce the quality of feature selection by optimizing for fewer category
 542 distinctions, it eliminated the possibility that mixing predictive and predicted categories may artificially in-
 543 flate classification performance. This approach again produced qualitatively similar results, though slightly
 544 weaker. We found reliable evidence for both A and B Current items, during the ON (A: mean = 0.0093; 95%
 545 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 =
 546 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
 547 was reliable during the ISI (mean = 0.0078; 95% CI = [0.0045, 0.012], $p < 0.001$), but not ON periods (mean =

548 0.0046; 95% CI = [-0.0012, 0.012], $p = 0.17$). Further, we numerically replicated our key finding of predictive
 549 B evidence during the Pre-ISI period (mean = 0.0038; 95% CI = [0.00, 0.0080], $p = 0.050$), though lingering A
 550 evidence during the Post-ON period was no longer reliable (mean = 0.0022; 95% CI = [-0.0034, 0.0081], $p =$
 551 0.47).

552 Taken together, these results show that statistical learning of the category pairs in Structured blocks
 553 affected neural representations in the task. Not only did visual contacts represent the category of the first
 554 and second items in a pair while they were being perceived (A and B evidence during ON and ISI epochs of
 555 A and B, respectively), but also the first category during the second (A evidence during ON epoch of B) and
 556 the second category during the first (B evidence during ISI epoch after A). This latter effect indicates that
 557 the first item in a pair (from A category) had predictive value on average.

558 We again examined whether these predictive effects emerged over time, in the first run of the Structured
 559 condition. For each participant, we computed the Spearman rank correlation of subblock number with
 560 the mean predictive evidence for B (averaged across all A items in each subblock), expecting a positive
 561 correlation. The resulting within-participant relationship was not reliable at the group level (mean rho =
 562 0.012; 95% CI = [-0.24, 0.24], $p = 0.92$). We also tested for a positive relationship across subblocks between
 563 prediction of B during A and neural entrainment for pairs, given that we expect both measures to depend
 564 upon statistical learning. However, this within-participant relationship was not reliable at the group level
 565 (mean rho = 0.038; 95% CI = [-0.12, 0.19], $p = 0.67$); nor was it reliable for neural entrainment to images
 566 (mean rho = -0.11; 95% CI = [-0.29, 0.079], $p = 0.25$).

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

569 Subsequent memory analysis

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

578 Consistent with our hypothesis, B evidence during the ISI epoch after A (i.e., Predicted category) was
 579 negatively related to subsequent A memory (**Figure 8A**): forgotten A items yielded reliable B evidence (mean
 580 = 0.0092; 95% CI = [0.0023, 0.017], $p = 0.0030$), whereas remembered A items did not (mean = 0.0017; 95%
 581 CI = [-0.0016, 0.0049], $p = 0.31$). In contrast, A evidence during the ISI epoch after A (i.e., Perceived category)
 582 was reliable for both remembered (mean = 0.012; 95% CI = [0.0091, 0.015], $p < 0.001$) and forgotten (mean =
 583 0.014; 95% CI = [0.0077, 0.021], $p < 0.001$) A items. This differential effect of subsequent memory on neural
 584 evidence for Perceived vs. Predicted categories during the ISI after A was reflected in a significant 2 (evidence
 585 category: A, B) by 2 (subsequent memory: remembered, forgotten) interaction ($p < 0.001$). This interaction
 586 was driven by a marginal difference in neural evidence for the Predicted B category during encoding of
 587 subsequently forgotten vs. remembered A items (mean difference = 0.0075; 95% CI = [-0.00046, 0.016], p
 588 = 0.065), but no reliable difference in neural evidence for the Perceived A category by subsequent memory
 589 (mean difference = 0.0022; 95% CI = [-0.0050, 0.0094], $p = 0.57$).

590 As a control analysis, we performed the key steps above in the Random blocks. These blocks did not
 591 contain pairs, and so we dummy-coded pairs of X items (X_1-X_2 instead of A-B). In contrast to Structured
 592 blocks, we did not expect that neural evidence of the "Predicted" X_2 category during the X_1 ISI would relate
 593 to subsequent memory for X_1 . Indeed, there was no reliable evidence for the X_2 category for either remem-
 594 bered (mean = -0.0029; 95% CI = [-0.0069, 0.00084], $p = 0.14$) or forgotten (mean = 0.0011; 95% CI = [-0.0027,
 595 0.0054], $p = 0.57$) X_1 items. In contrast, neural evidence for the Perceived X_1 category during the X_1 ISI was
 596 reliable for both remembered X_1 items (mean = 0.010; 95% CI = [0.0039, 0.019], $p < 0.001$) and forgotten X_1
 597 items (mean = 0.0065; 95% CI = [0.0022, 0.012], $p < 0.001$).

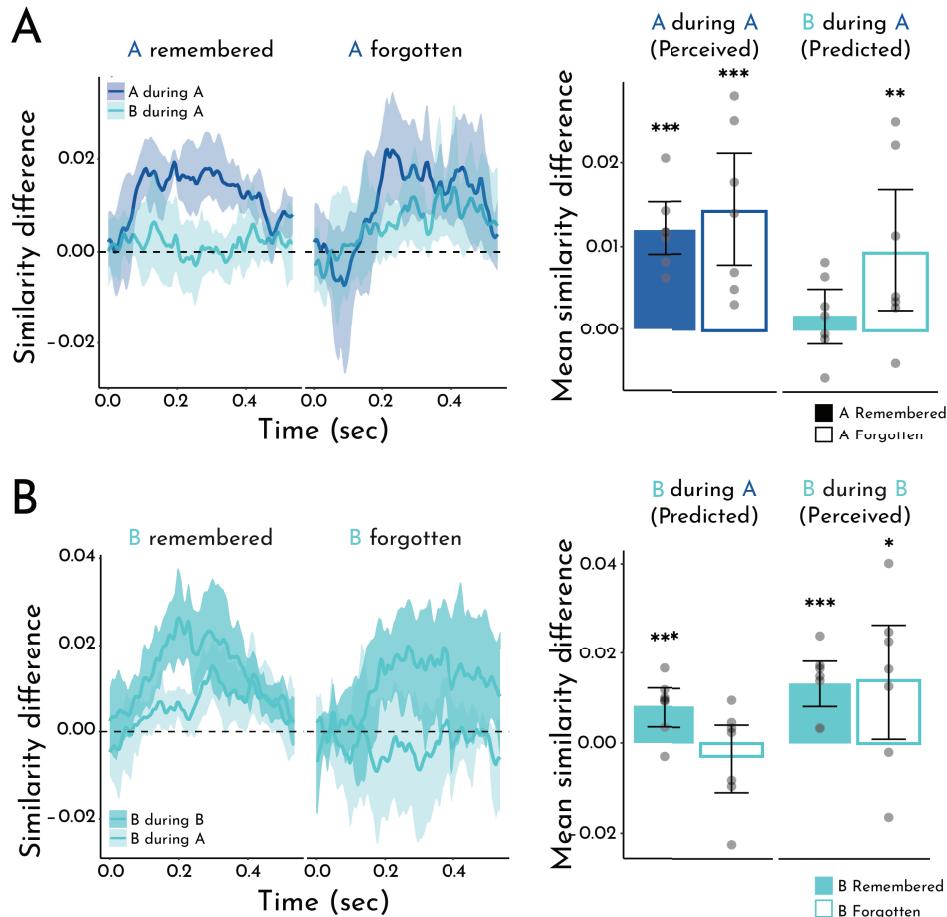


Figure 8. Subsequent memory analysis. A) Left: Timecourse of pattern similarity in visual contacts between A items being encoded and the Perceived A (A during A) and Predicted B (B during A) category templates, as a function of whether A items were subsequently remembered or forgotten. Right: Pattern 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 = remembered; empty bars = forgotten). B) Left: Timecourse of pattern similarity in visual contacts between B items being encoded and the Predicted B (B during A) and Perceived B (B during B) category templates, as a function of whether B items were subsequently remembered for forgotten. Right: Pattern similarity averaged within the ISI period, as a function of subsequent memory for B items. Error shading/bars reflect the bootstrapped 95% confidence interval across participants. Each dot represents an individual participant. * $p < 0.05$; ** $p < 0.01$; *** $p < 0.001$

598 We so far focused on the effects of prediction for memory of the item generating the prediction (A), but
599 what is the mnemonic fate of the item being predicted (B), which in this task with deterministic pairs always
600 appeared as expected? Whereas neural category evidence for B during the A ISI (Predicted) was negatively
601 related to subsequent memory for A items, the opposite was true for memory of B items (**Figure 8B**): re-
602 membered B items were associated with reliable prediction of B (mean = 0.0082; 95% CI = [0.0036, 0.012], p
603 <0.001), but forgotten B items were not (mean = -0.0028; 95% CI = [-0.011, 0.0041], p = 0.49). In contrast, and
604 similar to A memory, evidence for B during the B ISI (Perceived) was reliable for both remembered (mean =
605 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)
606 B items. We did not find an interaction between category and memory (p = 0.22). However, there was a
607 reliable difference in Predicted B evidence for remembered vs. forgotten B items (mean difference = 0.011;
608 95% CI = [0.00060, 0.021], p = 0.039); Perceived B evidence did not differ as a function of memory (mean
609 difference = 0.00064; 95% CI = [-0.014, 0.016], p = 0.89).

610 We repeated the same control analysis of Random blocks, but now focused on subsequent memory for
611 X_2 items (equivalent to B, rather than X_1 memory for A). Neural evidence for the "Predicted" X_2 category
612 during the ISI after X_1 was not reliable for either remembered (mean = 0.0013; 95% CI = [-0.0020, 0.0043], p
613 = 0.44) or forgotten (mean = -0.00048; 95% CI = [-0.0030, 0.0017], p = 0.75) X_2 items.

614 We again tested whether our key results generalized to templates created from two alternative classifica-
615 tion approaches. Using a 6-way classifier, we replicated the finding that forgotten A items were associated
616 with reliable predictive evidence of B (mean = 0.0075; 95% CI = [0.0015, 0.014], p = 0.009), whereas remem-
617 bered A items were not (mean = 0.0026; 95% CI = [-0.00010, 0.0054], p = 0.061). In contrast, forgotten B
618 items were not associated with reliable predictive evidence of B (mean = -0.0046; 95% CI = [-0.016, 0.0037],
619 p = 0.40), whereas remembered B items were (mean = 0.0082; 95% CI = [0.0021, 0.015], p = 0.003). Using
620 binary classifiers trained to discriminate within A or B categories, we again found that forgotten (mean =
621 0.0075; 95% CI = [0.00087, 0.016], p = 0.014), but not remembered A items (mean = 0.0027; 95% CI = [-
622 0.00086, 0.0061], p = 0.13) were associated with reliable predictive evidence of B, and that remembered
623 (mean = 0.0084; 95% CI = [0.0033, 0.013], p = 0.0016), but not forgotten B items (mean = -0.0044; 95% CI =
624 [-0.017, 0.0048], p = 0.47) were associated with reliable predictive evidence of B.

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

630 Discussion

631 This study demonstrates a trade-off between how well an item is encoded into episodic memory and how
632 strong of a future prediction it generates based on statistical learning. We first used frequency tagging
633 to provide neural verification of statistical learning. During a sequence of scene photographs, electrodes
634 in visual cortex represented pairs of scene categories that reliably followed each other, synchronizing not
635 only to the individual scenes but also to the boundaries between pairs. Next, we used multivariate pattern
636 analysis to assess how the paired categories were represented over time. Items from the first category in a
637 pair elicited a representation of the second category, which grew in strength in advance of the onset of items
638 from the second category. We refer to the ability of an item to generate this predictive representation as its
639 "predictive value". Critically, by relating these representational dynamics to subsequent memory behavior,
640 we found that forgotten items from the first category triggered reliable predictions during encoding whereas
641 remembered first items had not.

642 Our work builds upon suggestive evidence from a prior study that predictive value may influence sub-
643 sequent memory (**Sherman and Turk-Browne, 2020**). This prior study included behavioral and fMRI experi-
644 ments, whereas the current study employed iEEG. Neural measures are an important advance over behav-
645 ior alone because they can assay predictive representations during passive viewing at encoding. iEEG is
646 superior to fMRI for this purpose because neural activity is sampled at much greater temporal resolution
647 and activity reflects instantaneous electrical potentials rather than hemodynamic responses smoothed and

648 delayed in time. This provides much greater confidence that the upcoming category was being represented
649 *prior* to its appearance and thus was truly predictive. Moreover, the prior study showed a negative rela-
650 tionship between prediction and memory across participants, whereas the current study established this
651 relationship within participant. This is also an important advance because an across-participant relation-
652 ship does not provide strong evidence for the claim that prediction during encoding impairs memory. Such
653 a relationship could reflect generic individual differences such that, for example, a participant with better
654 overall memory generates the same weak prediction on both remembered and forgotten trials. In contrast,
655 in this study we were able to link prediction to successful vs. unsuccessful memory formation across items.
656 This more sensitive approach yielded other findings not observed in the prior study, including that memory
657 for B items had an opposite, positive relationship with prediction of B. Taken together, these results pro-
658 vide mechanistic insight into the interaction between predictive value and memory, and speak to theoretical
659 questions about the representations underlying statistical learning and episodic memory.

660 **Nature of representational changes**

661 Several fMRI studies have shown that statistical and related forms of learning can change neural representa-
662 tions of associated items throughout the human brain (**Schapiro et al., 2012, 2013; Schlichting et al., 2015;**
663 **Deuker et al., 2016; Tompary and Davachi, 2017**). For example, if exposed to sequential pairs embedded
664 in a continuous stream of objects (akin to the category pairs in the current study), the two objects in a pair
665 come to elicit more similar patterns of fMRI activity from before to after learning, when presented on their
666 own, in the medial temporal lobe cortex and hippocampus (**Schapiro et al., 2012**). Such integration could
667 be interpreted as evidence that the representations of the paired items merged into a single “unitized” rep-
668 resentation of the pair that can be evoked by either item (**Fujimichi et al., 2010**). Alternatively, the paired
669 items may remain distinct but become associated, such that either can be reactivated by the other through
670 spreading activation (**Schapiro et al., 2017**). A key difference between these accounts is the timing of how
671 learned representations emerge when one of the items is presented: the merging account predicts that
672 the (same) unitized representation is evoked immediately by either paired item, whereas the associative
673 account predicts that the presented item is represented immediately while the paired item is represented
674 gradually over time through reactivation. These dynamics cannot be distinguished by fMRI because of its
675 slow temporal resolution, but our iEEG approach may shed light.

676 On the surface, the results of our frequency tagging analysis may seem to suggest a merged represen-
677 tation of the category pairs. The reliable peak in coherence at the frequency of two consecutive stimuli may
678 suggest that electrodes in visual cortex represented the paired categories as a single unit (**Batterink and Paller,**
679 **2017**). However, the results of our pattern similarity analysis are more consistent with an association be-
680 tween the paired categories. Although we found that both categories in a pair could be represented at
681 the same time (i.e., predictive B evidence during the A Pre trial and lingering A evidence during the B Post
682 trial, relative to no such evidence on X trials), these representations were offset in time. The representation
683 of the A category was robust during both the ON and ISI epochs of the A trial, whereas the representa-
684 tion of the B category was not reliable during the ON epoch and only emerged during the ISI epoch. Thus,
685 our results are more consistent with an associative account in visual cortex. It remains possible that the
686 hippocampus or other brain structures represent statistical regularities through unitized representations.
687 Moreover, one limitation of our study is that we did not measure representations of individual categories be-
688 fore and after learning to directly assess representational change. Although we could not directly measure
689 representational change from before to after learning, we did correlate the category templates measured
690 after learning. Unitization of paired categories would be reflected in increased pattern similarity among
691 paired, relative to unpaired and random categories. We did not find reliable evidence of such representa-
692 tional merging, inconsistent with a unitization account. However, prior studies focused on the unitization of
693 paired items rather than categories. Thus, if we had found evidence of representational merging of paired
694 categories in the current study, it would be unclear whether this reflects unitization in the same way or a
695 qualitatively different kind of representational change.

696 Predictive interference on memory encoding

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

707 Interestingly, evidence for the current A category was comparable across remembered and forgotten
708 A items. Thus, in this paradigm, variance in memory was explained solely by prediction of the upcoming
709 category, not the strength of perceptual processing of the category being encoded (*Kuhl et al., 2012*) nor
710 modulation of this processing by prediction (both of which would have affected A evidence). The lack of a
711 relationship between A evidence and A memory may reflect a tradeoff: category evidence may reflect rep-
712 resentation of the most diagnostic features of a category, which would enhance memory for these features
713 while impairing memory for idiosyncratic features of particular exemplars. A related account may explain
714 why predictive B evidence was positively linked to B memory (*Smith et al., 2013; Thavabalasingam et al.,*
715 **2016**): B evidence during the A ISI may potentiate the diagnostic features of the B category, enhancing the
716 salience of idiosyncratic features of B when it appears to strengthen episodic memory for B. Future studies
717 could test these possibilities by using a more continuous measure of memory precision and by testing on
718 modified items that retain category-diagnostic vs. idiosyncratic features.

719 Our finding that prediction relates to better memory for predictable B items contrasts with findings
720 of enhanced encoding for unpredictable/unexpected items (*Kim et al., 2014; Greve et al., 2017; Bein et al.,*
721 **2021**). These seemingly divergent findings are difficult to reconcile because predictions in our study were
722 never violated: in the Structured condition, the A in each pair was followed deterministically by B; in the
723 Random condition, although each X was unexpected to some degree they did not violate a learned expecta-
724 tion. Thus, it is possible that replacing the expected B with another category would have led to even better
725 memory encoding. That said, one interpretation of our finding of enhanced (predictable) B memory that
726 would be consistent with a benefit of prediction error for episodic memory could be that features idiosyn-
727 cratic to a particular B exemplar (needed to later retrieve this specific episodic memory) may have violated a
728 category-level expectation grounded in the diagnostic (i.e., non-idiosyncratic) features of a category shared
729 across its exemplars. This question — as well as questions above about how the category-level nature of
730 the prediction may have affected memory for A — could be informed by future studies examining effects
731 of item-level prediction on memory.

732 This work builds on existing theories considering the complex interplay between memory encoding and
733 memory retrieval. To the extent that prediction from statistical learning can be considered associative
734 retrieval (*Kok and Turk-Browne, 2018; Hindy et al., 2016*), our findings converge with the notion that the
735 brain cycles between mutually exclusive encoding and retrieval states (*Hasselmo et al., 2002; Duncan et al.,*
736 **2012; Long and Kuhl, 2019; Bein et al., 2020**), organized by the hippocampal theta cycle (*Kerrén et al., 2018;*
737 *Pacheco Estefan et al., 2021*). Further, a recent computational model suggests that predictive uncertainty
738 determines when memories should be encoded or retrieved (*Lu et al., 2022*). The model accounts for find-
739 ings that familiar experiences are more likely to evoke retrieval (*Patil and Duncan, 2018*), and thus may help
740 to explain why predictions from statistical learning are prioritized over episodic encoding.

741 Neural source of predictions

742 The current study sought to decode evidence of visual categories and so focused on electrode contacts in
743 visual cortex. This adds to a growing literature on predictive signals in visual cortex (*De Lange et al., 2018;*
744 *Kim et al., 2020; Clarke et al., 2022*). Importantly, in our previous fMRI study (*Sherman and Turk-Browne,*
745 **2020**), we found evidence of prediction only in the hippocampus. We interpreted the lack of an effect in

746 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 on-screen, 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 vs. ISI time periods when there was vs. was not a competing stimulus present, respectively (which fMRI would have been unable to separate). In fact, 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.

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

762 Beyond generating predictions, the hippocampus may also be the nexus of the interaction between
763 episodic memory and statistical learning, given its fundamental role in both functions (**Schapiro et al., 2017**).
764 Indeed, given the necessity of the hippocampus for episodic memory, our study raises questions about how
765 the representations of perceived and predicted categories in visual cortex are routed into the hippocampus
766 for encoding. One intriguing possibility is that these representations are prioritized according to biased competition
767 (**Desimone, 1998; Hutchinson et al., 2016**), leading to preferential routing and subsequent encoding
768 of predicted, but not perceived, information in the hippocampus. Relatedly, recent work had found that encoding
769 vs. retrieval states are associated with distinct patterns of activity in visual cortex (**Long and Kuhl,**
770 **2021**), suggesting that representations in visual regions may be fundamentally shaped by memory state in
771 the hippocampus.

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

780 **Limitations of the current study**

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

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

796 2013).

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

814 Lastly, we adopted a “subblock” structure, in which individual exemplars repeated four times before
815 switching to new exemplars (but holding the category pairs constant). This choice was made to balance
816 the rapid presentation of stimuli needed for the neural frequency tagging analyses with providing sufficient
817 exposure to the images so that some would be later remembered. Although we found some evidence that
818 neural entrainment to the pairs increased across Structured subblocks, there was little evidence of a learn-
819 ing trajectory in the behavioral or predictive neural measures. It is possible that exemplar repetition in the
820 subblocks may have allowed learning to asymptote after only one or a few subblocks (**Turk-Browne et al.,**
821 **2009**), eliminating the possibility of finding a more gradual change in these measures across subblocks.
822 These analyses are further limited by the small number of patients relative to prior work with healthy in-
823 dividuals that found clearer learning effects in behavior (**Sherman and Turk-Browne, 2020**). Future studies
824 could tailor their experimental designs to optimize detection of a learning trajectory, for example by forego-
825 ing neural entrainment and presenting images once for a longer duration or by introducing more complex
826 regularities.

827 Conclusion

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

834 Acknowledgments

835 We are grateful to the patients who participated in this study. We thank Kun Wu for providing the elec-
836 trode reconstructions, Christopher Benjamin for helping to recruit patients and coordinate testing, Richard
837 Aslin and Sami Yousif for helpful conversations, and Gregory McCarthy for advice about data collection and
838 analysis, as well as for feedback on the manuscript. This work was supported by NIH grant R01 MH069456
839 (N.B.T-B.), the Canadian Institute for Advanced Research (N.B.T-B.), and an NSF GRFP grant (B.E.S.).

840 Competing Interests

841 The authors declare no competing interests.

842 References

- 843 Aitken F, Kok P. Hippocampal representations switch from errors to predictions during acquisition of predictive associations. *Nature Communications*. 2022; 13(1):1–13.
- 845 Aitken F, Turner G, Kok P. Prior expectations of motion direction modulate early sensory processing. *Journal of Neuroscience*. 2020; 40(33):6389–6397.
- 847 Aly M, Turk-Browne NB. How hippocampal memory shapes, and is shaped by, attention. In: *The hippocampus from cells to systems* Springer; 2017.p. 369–403.
- 849 Batterink LJ, Paller KA. Online neural monitoring of statistical learning. *Cortex*. 2017; 90:31–45.
- 850 Bein O, Duncan K, Davachi L. Mnemonic prediction errors bias hippocampal states. *Nature communications*. 2020; 11(1):1–11.
- 852 Bein O, Plotkin NA, Davachi L. Mnemonic prediction errors promote detailed memories. *Learning and Memory*. 2021; .
- 853 Biderman N, Bakkour A, Shohamy D. What are memories for? The hippocampus bridges past experience with future decisions. *Trends in Cognitive Sciences*. 2020; 24(7):542–556.
- 855 Bosch SE, Jehee JF, Fernández G, Doeller CF. Reinstatement of associative memories in early visual cortex is signaled by the hippocampus. *Journal of Neuroscience*. 2014; 34(22):7493–7500.
- 857 Brady TF, Oliva A. Statistical learning using real-world scenes: Extracting categorical regularities without conscious intent. *Psychological science*. 2008; 19(7):678–685.
- 859 Choi D, Batterink LJ, Black AK, Paller KA, Werker JF. Preverbal Infants Discover Statistical Word Patterns at Similar Rates as Adults: Evidence From Neural Entrainment. *Psychological Science*. 2020; 31(9):1161–1173.
- 861 Clarke A, Crivelli-Decker J, Ranganath C. Contextual expectations shape cortical reinstatement of sensory representations. *Journal of Neuroscience*. 2022; 42(30):5956–5965.
- 863 Cowan ET, Schapiro AC, Dunsmoor JE, Murty VP. Memory consolidation as an adaptive process. *Psychonomic Bulletin & Review*. 2021; 28(6):1796–1810.
- 865 Danker JF, Tompary A, Davachi L. Trial-by-trial hippocampal encoding activation predicts the fidelity of cortical reinstatement during subsequent retrieval. *Cerebral Cortex*. 2017; 27(7):3515–3524.
- 867 De Brigard F. Is memory for remembering? Recollection as a form of episodic hypothetical thinking. *Synthese*. 2014; 191(2):155–185.
- 869 De Lange FP, Heilbron M, Kok P. How do expectations shape perception? *Trends in cognitive sciences*. 2018; 22(9):764–779.
- 871 Demarchi G, Sanchez G, Weisz N. Automatic and feature-specific prediction-related neural activity in the human auditory system. *Nature communications*. 2019; 10(1):1–11.
- 873 Desimone R. Visual attention mediated by biased competition in extrastriate visual cortex. *Philosophical Transactions of the Royal Society of London Series B: Biological Sciences*. 1998; 353(1373):1245–1255.
- 875 Deuker L, Bellmund JL, Schröder TN, Doeller CF. An event map of memory space in the hippocampus. *Elife*. 2016; 5:e16534.
- 877 Dickerson KC, Adcock RA. Motivation and memory. *Stevens' Handbook of Experimental Psychology and Cognitive Neuroscience, Learning and Memory*. 2018; 1:215.
- 879 Dolcos F, Katsumi Y, Weymar M, Moore M, Tsukiura T, Dolcos S. Emerging directions in emotional episodic memory. *Frontiers in Psychology*. 2017; 8:1867.
- 881 Duncan K, Sadanand A, Davachi L. Memory's penumbra: episodic memory decisions induce lingering mnemonic biases. *Science*. 2012; 337(6093):485–487.
- 883 Efron B, Tibshirani R. Bootstrap methods for standard errors, confidence intervals, and other measures of statistical accuracy. *Statistical science*. 1986; p. 54–75.
- 885 Endress AD, Johnson SP. When forgetting fosters learning: A neural network model for statistical learning. *Cognition*. 2021; p. 104621.

- 887 Fiser J, Aslin RN. Statistical learning of higher-order temporal structure from visual shape sequences. *Journal of Experimental Psychology: Learning, Memory, and Cognition*. 2002; 28(3):458.
- 888
- 889 Fujimichi R, Naya Y, Koyano KW, Takeda M, Takeuchi D, Miyashita Y. Unitized representation of paired objects in area
890 35 of the macaque perirhinal cortex. *European Journal of Neuroscience*. 2010; 32(4):659–667.
- 891 Gebhart AL, Aslin RN, Newport EL. Changing structures in midstream: Learning along the statistical garden path. *Cognitive science*. 2009; 33(6):1087–1116.
- 892
- 893 Goldfarb EV. Enhancing memory with stress: progress, challenges, and opportunities. *Brain and Cognition*. 2019; 133:94–
894 105.
- 895 Gómez DM, Bion RA, Mehler J. The word segmentation process as revealed by click detection. *Language and Cognitive Processes*. 2011; 26(2):212–223.
- 896
- 897 Greve A, Cooper E, Kaula A, Anderson MC, Henson R. Does prediction error drive one-shot declarative learning? *Journal of memory and language*. 2017; 94:149–165.
- 898
- 899 Grier JB. Nonparametric indexes for sensitivity and bias: computing formulas. *Psychological bulletin*. 1971; 75(6):424.
- 900 Hasselmo ME, Bodelón C, Wyble BP. A proposed function for hippocampal theta rhythm: separate phases of encoding
901 and retrieval enhance reversal of prior learning. *Neural computation*. 2002; 14(4):793–817.
- 902 Henin S, Turk-Browne NB, Friedman D, Liu A, Dugan P, Flinker A, Doyle W, Devinsky O, Melloni L. Learning hierarchical
903 sequence representations across human cortex and hippocampus. *Science advances*. 2021; 7(8):eabc4530.
- 904 Hindy NC, Ng FY, Turk-Browne NB. Linking pattern completion in the hippocampus to predictive coding in visual cortex.
905 *Nature neuroscience*. 2016; 19(5):665–667.
- 906 Hutchinson JB, Pak SS, Turk-Browne NB. Biased competition during long-term memory formation. *Journal of cognitive
907 neuroscience*. 2016; 28(1):187–197.
- 908 Jenkinson M, Bannister P, Brady M, Smith S. Improved optimization for the robust and accurate linear registration and
909 motion correction of brain images. *Neuroimage*. 2002; 17(2):825–841.
- 910 Jenkinson M, Beckmann CF, Behrens TE, Woolrich MW, Smith SM. Fsl. *Neuroimage*. 2012; 62(2):782–790.
- 911 Jenkinson M, Smith S. A global optimisation method for robust affine registration of brain images. *Medical image
912 analysis*. 2001; 5(2):143–156.
- 913 Jungé JA, Scholl Bj, Chun MM. How is spatial context learning integrated over signal versus noise? A primacy effect in
914 contextual cueing. *Visual cognition*. 2007; 15(1):1–11.
- 915 Karuza EA, Newport EL, Aslin RN, Starling SJ, Tivarus ME, Bavelier D. The neural correlates of statistical learning in a word
916 segmentation task: An fMRI study. *Brain and language*. 2013; 127(1):46–54.
- 917 Kerrén C, Linde-Domingo J, Hanslmayr S, Wimber M. An optimal oscillatory phase for pattern reactivation during memory
918 retrieval. *Current Biology*. 2018; 28(21):3383–3392.
- 919 Kim G, Lewis-Peacock JA, Norman KA, Turk-Browne NB. Pruning of memories by context-based prediction error. *Proceedings of the National Academy of Sciences*. 2014; 111(24):8997–9002.
- 920
- 921 Kim H, Schlichting ML, Preston AR, Lewis-Peacock JA. Predictability changes what we remember in familiar temporal
922 contexts. *Journal of cognitive neuroscience*. 2020; 32(1):124–140.
- 923 Kok P, Failing MF, de Lange FP. Prior expectations evoke stimulus templates in the primary visual cortex. *Journal of cognitive
924 neuroscience*. 2014; 26(7):1546–1554.
- 925 Kok P, Mostert P, De Lange FP. Prior expectations induce prestimulus sensory templates. *Proceedings of the National
926 Academy of Sciences*. 2017; 114(39):10473–10478.
- 927 Kok P, Turk-Browne NB. Associative prediction of visual shape in the hippocampus. *Journal of Neuroscience*. 2018;
928 38(31):6888–6899.
- 929 Kuhl BA, Rissman J, Wagner AD. Multi-voxel patterns of visual category representation during episodic encoding are
930 predictive of subsequent memory. *Neuropsychologia*. 2012; 50(4):458–469.

- 931 Long NM, Kuhl BA. Decoding the tradeoff between encoding and retrieval to predict memory for overlapping events.
932 *NeuroImage*. 2019; 201:116001.
- 933 Long NM, Kuhl BA. Cortical representations of visual stimuli shift locations with changes in memory states. *Current*
934 *Biology*. 2021; 31(5):1119–1126.
- 935 Lu Q, Hasson U, Norman KA. A neural network model of when to retrieve and encode episodic memories. *eLife*. 2022;
936 11:e74445.
- 937 Oostenveld R, Fries P, Maris E, Schoffelen JM. FieldTrip: open source software for advanced analysis of MEG, EEG, and
938 invasive electrophysiological data. *Computational intelligence and neuroscience*. 2011; 2011.
- 939 Pacheco Estefan D, Zucca R, Arsiwalla X, Principe A, Zhang H, Rocamora R, Axmacher N, Verschure PF. Volitional learning
940 promotes theta phase coding in the human hippocampus. *Proceedings of the National Academy of Sciences*. 2021;
941 118(10):e2021238118.
- 942 Papademetris X, Jackowski MP, Rajeevan N, DiStasio M, Okuda H, Constable RT, Staib LH. BiolImage Suite: An integrated
943 medical image analysis suite: An update. *The insight journal*. 2006; 2006:209.
- 944 Patil A, Duncan K. Lingering cognitive states shape fundamental mnemonic abilities. *Psychological Science*. 2018;
945 29(1):45–55.
- 946 Reddy L, Self MW, Zoefel B, Poncet M, Posse JK, Peters JC, Baayen JC, Idema S, VanRullen R, Roelfsema PR. Theta-phase de-
947 pendent neuronal coding during sequence learning in human single neurons. *Nature communications*. 2021; 12(1):1–9.
- 948 Schacter DL, Addis DR, Hassabis D, Martin VC, Spreng RN, Szpunar KK. The future of memory: remembering, imagining,
949 and the brain. *Neuron*. 2012; 76(4):677–694.
- 950 Schapiro AC, Kustner LV, Turk-Browne NB. Shaping of object representations in the human medial temporal lobe based
951 on temporal regularities. *Current biology*. 2012; 22(17):1622–1627.
- 952 Schapiro AC, Rogers TT, Cordova NI, Turk-Browne NB, Botvinick MM. Neural representations of events arise from tem-
953 poral community structure. *Nature neuroscience*. 2013; 16(4):486–492.
- 954 Schapiro AC, Turk-Browne NB, Botvinick MM, Norman KA. Complementary learning systems within the hippocampus:
955 a neural network modelling approach to reconciling episodic memory with statistical learning. *Philosophical Transac-
956 tions of the Royal Society B: Biological Sciences*. 2017; 372(1711):20160049.
- 957 Schlichting ML, Mumford JA, Preston AR. Learning-related representational changes reveal dissociable integration and
958 separation signatures in the hippocampus and prefrontal cortex. *Nature communications*. 2015; 6(1):1–10.
- 959 Sherman BE, Graves KN, Turk-Browne NB. The prevalence and importance of statistical learning in human cognition and
960 behavior. *Current opinion in behavioral sciences*. 2020; 32:15–20.
- 961 Sherman BE, Turk-Browne NB. Statistical prediction of the future impairs episodic encoding of the present. *Proceedings*
962 *of the National Academy of Sciences*. 2020; 117(37):22760–22770.
- 963 Siegelman N, Bogaerts L, Kronenfeld O, Frost R. Redefining “learning” in statistical learning: What does an online measure
964 reveal about the assimilation of visual regularities? *Cognitive science*. 2018; 42:692–727.
- 965 Smith TA, Hasinski AE, Sederberg PB. The context repetition effect: Predicted events are remembered better, even when
966 they don’t happen. *Journal of Experimental Psychology: General*. 2013; 142(4):1298.
- 967 Tanaka KZ, Pevzner A, Hamidi AB, Nakazawa Y, Graham J, Wiltgen BJ. Cortical representations are reinstated by the
968 hippocampus during memory retrieval. *Neuron*. 2014; 84(2):347–354.
- 969 Thavabalasingam S, O’Neil EB, Zeng Z, Lee AC. Recognition memory is improved by a structured temporal framework
970 during encoding. *Frontiers in Psychology*. 2016; 6:2062.
- 971 Tompary A, Davachi L. Consolidation promotes the emergence of representational overlap in the hippocampus and
972 medial prefrontal cortex. *Neuron*. 2017; 96(1):228–241.
- 973 Treder MS, Charest I, Michelmann S, Martín-Buro MC, Roux F, Carceller-Benito F, Ugaldé-Canitrot A, Rollings DT, Sawlani
974 V, Chevallierah R, et al. The hippocampus as the switchboard between perception and memory. *Proceedings of the*
975 *National Academy of Sciences*. 2021; 118(50).
- 976 Turk-Browne NB, Jungé JA, Scholl BJ. The automaticity of visual statistical learning. *Journal of Experimental Psychology:*
977 *General*. 2005; 134(4):552.

- 978 Turk-Browne NB, Scholl BJ, Chun MM, Johnson MK. Neural evidence of statistical learning: Efficient detection of visual
979 regularities without awareness. *Journal of cognitive neuroscience*. 2009; 21(10):1934–1945.
- 980 Walther DB, Caddigan E, Fei-Fei L, Beck DM. Natural scene categories revealed in distributed patterns of activity in the
981 human brain. *Journal of neuroscience*. 2009; 29(34):10573–10581.