

# Long-Horizon Associative Learning Explains Human Sensitivity to Statistical and Network Structures in Auditory Sequences

Lucas Benjamin,<sup>1</sup> Mathias Sablé-Meyer,<sup>1,2</sup> Ana Fló,<sup>1,3</sup> Ghislaine Dehaene-Lambertz,<sup>1\*</sup> and Fosca Al Roumi<sup>1\*</sup>

<sup>1</sup>Cognitive Neuroimaging Unit, CNRS ERL 9003, INSERM U992, CEA, Université Paris-Saclay, NeuroSpin Center, 91190 Gif/Yvette, France, <sup>2</sup>Sainsbury Wellcome Centre for Neural Circuits and Behaviour, University College London, London W1T 4JG, United Kingdom, and <sup>3</sup>Department of Developmental Psychology and Socialization, University of Padova, Padova 35131, Italy

Networks are a useful mathematical tool for capturing the complexity of the world. In a previous behavioral study, we showed that human adults were sensitive to the high-level network structure underlying auditory sequences, even when presented with incomplete information. Their performance was best explained by a mathematical model compatible with associative learning principles, based on the integration of the transition probabilities between adjacent and nonadjacent elements with a memory decay. In the present study, we explored the neural correlates of this hypothesis via magnetoencephalography (MEG). Participants ( $N = 23$ , 16 females) passively listened to sequences of tones organized in a sparse community network structure comprising two communities. An early difference ( $\sim 150$  ms) was observed in the brain responses to tone transitions with similar transition probability but occurring either within or between communities. This result implies a rapid and automatic encoding of the sequence structure. Using time-resolved decoding, we estimated the duration and overlap of the representation of each tone. The decoding performance exhibited exponential decay, resulting in a significant overlap between the representations of successive tones. Based on this extended decay profile, we estimated a long-horizon associative learning novelty index for each transition and found a correlation of this measure with the MEG signal. Overall, our study sheds light on the neural mechanisms underlying human sensitivity to network structures and highlights the potential role of Hebbian-like mechanisms in supporting learning at various temporal scales.

**Key words:** associative learning; community structure; MEG; network learning; sequence; statistical learning

## Significance Statement

We conducted a MEG study in which human adults were passively exposed to sequences of tones organized in a sparse community network structure. Despite the uniform transition probabilities between tones, participants' brain activity exhibited sensitivity to the network structure. Notably, a consistent “deviant” response was observed at  $\sim 150$  ms when the sequence switched between communities. A long-tail exponential decay in tone representation allowed overlapping representations of successive sequence elements, facilitating long-range associative mechanisms. This binding mechanism adequately accounted for various scales of sequence learning, bridging the gap between statistical and network learning approaches.

## Introduction

Understanding the structure of the input sequences we encounter is fundamental for developing a comprehensive mental model of our environment (Dehaene et al., 2015, 2022). The capacity to

detect first-order relationships between successive events (i.e., transition probabilities) and its limits have been extensively studied in humans at the behavioral and neural levels (Saffran et al., 1996; Maheu et al., 2019; Benjamin et al., 2021, 2023b, 2024; Henin et al., 2021; Fló et al., 2022) as well as in nonhumans animals (Toro and Trobalón, 2005; James et al., 2020; Boros et al., 2021). Higher-order statistical relations between elements of a sequence are also detected by human adults and children (Schapiro et al., 2013; Karuza et al., 2019; Lynn et al., 2020; Mark et al., 2020), but only a limited number of neuroimaging studies have explored possible neural correlates of this learning (Schapiro et al., 2016; Ren et al., 2022; Stiso et al., 2022). Therefore, we still do not know whether a single mechanism

Received July 4, 2023; revised Jan. 16, 2024; accepted Feb. 7, 2024.

Author contributions: L.B., A.F., G.D.-L., and F.A.R. designed research; L.B. and M.S.-M. performed research; L.B., M.S.-M., and F.A.R. analyzed data; L.B., G.D.-L., and F.A.R. wrote the paper.

We thank Lucia Melloni for discussions and remarks during the data analysis.

\*G.D.-L. and F.A.R. contributed equally to this work.

The authors declare no competing financial interests.

Correspondence should be addressed to Lucas Benjamin at lucas.benjamin78@gmail.com.

<https://doi.org/10.1523/JNEUROSCI.1369-23.2024>

Copyright © 2024 the authors

can adequately explain both first-order (local transitions) and network structure learning and whether these computations require distinct cognitive and brain processes.

To bridge the gap between local statistical and network-level learning studies, we previously proposed the sparse community paradigm, which allows to simultaneously characterize these aspects on auditory sequences (Benjamin et al., 2023a). Building upon the community paradigm introduced by Schapiro et al. (2013), we created a network consisting of two densely but incompletely connected clusters (called communities) of six elements each. Each element is connected with four other elements, and the two clusters are linked by only two edges (links). A learning sequence is created by randomly drawing the next tone from the four possibilities, creating a random walk in the network with a uniform transition probability (TP) between successive tones (Movie 1A). After exposure to such a sequence, participants were asked to judge their familiarity with various pairs of tones that (1) had or had not been presented during learning to test local TP learning and (2) did or did not belong to the same community to test learning of the higher-level structure (Benjamin et al., 2023a). Interestingly, participants judged new transitions they had never heard as highly familiar if they were between tones belonging to the same community. This completion effect demonstrated that they generalized the community structure to missing transitions. Conversely, they judged transitions between communities to be less familiar than within communities despite the absence of any difference in local TP during learning. This pruning effect translates into a decrease in subjective familiarity with tone pairs that switch from one community to the other despite similar transition probabilities between tones. Among the various models proposed in the statistical and network learning literature, an associative learning approach (the free energy minimization model—FEMM; Lynn et al., 2020), conceptually related to the successor representation, provided the best fit to participants' behavior. According to this model, participants did not solely compute adjacent transition probabilities but a linear sum of transition probabilities at all orders (adjacent, first-order nonadjacent, second-order nonadjacent, and so on), weighted by a decreasing exponential factor. This model explains how both local transitions and network structures are perceived and successfully accounts for behavioral results across different network types, including community, sparse community, ring, and lattice networks (Lynn et al., 2020; Benjamin et al., 2023a), as well as results concerning local statistical learning. FEMM appears to be a good candidate for a unifying framework of sequence learning.

However, a common model is insufficient to postulate a common implementation (Marr, 1982), and there is still no consensus on how the brain implements these computations. On the one hand, sensitivity to network structure is often described as a conscious abstraction of the structure involving top-down attention processes with late brain signatures (Ren et al., 2022) typically in the prefrontal cortex (Stiso et al., 2022). On the other hand, we previously postulated that low-level associative learning (Benjamin et al., 2023a; see also Endress, 2010; Schapiro et al., 2017; Endress and Johnson, 2021) was sufficient for both local and higher-order learning. To disentangle those two hypotheses, we tested here whether passive exposure to a rapid auditory sequence could lead to successful learning of its network structure. We thus exposed participants to fast sequence of tones following the sparse community design while recording their brain activity with MEG.

## Materials and Methods

### Stimuli and procedure

We generated 12 tones of 50 ms duration, logarithmically distributed from 300 to 1,800 Hz. For each participant, the 12 tones were randomly assigned to the 12 nodes of the sparse community network (see Movie 1 for a complete description of the network structure). The sparse community network comprised two communities (i.e., clusters) made of six nodes, densely connected to each other but poorly connected to the nodes of the other community. Crucially, in the sparse community design, some connections between nodes belonging to the same community are missing. Specifically, for each participant, we randomly removed 12 transitions (six per community, one per node). After a training block with this incomplete graph, new transitions were added at a low frequency (4%) for the following test blocks. We refer to these transitions as New and those presented during training as Familiar. The New transitions were critically within and between the communities (which we refer to as Within vs Between). New Within corresponds to the 12 “missing” transitions randomly removed from the network in the training block. To balance, we randomly selected 12 New Between transitions (one per node) that violated the community clustering property. As a result, the transition probabilities between tones during the training block were flat: TP = 25%, while during the test blocks, the Familiar transitions had TP = 23% and a frequency of 18.4/block. The 12 New Within and 12 New Between community transitions had TP = 4% and a frequency of 3.2/block. The New Within and New Between transitions were randomly drawn for each subject to add variability to the network structure. Movie 1A shows an example of one structure and the associated sequence used for one participant.

We then performed random walks in the participant's sparse community graph to derive one 960-item-long training sequence and six 960-item-long test sequences (one sequence corresponds to one block) with 200 ms interstimulus interval (ISI) between each tone (Movie 1A). The first block of 960 items comprised only Familiar Within and Familiar Between transitions (training block, TP = 25% each). For the next six blocks (Tests 1–6), we introduced infrequent New Within and New Between transitions (TP = 4% each). All Familiar transitions, independently of whether they were Within or Between communities, had the same TPs and appeared with the same frequency (TP = 23% each). However, the graph structure entails that the participants heard in total fewer between community transitions than within community transitions (there are 32 Familiar Within and 4 Familiar Between community transitions during training, completed by 12 New Within and 12 New Between community transitions during test).

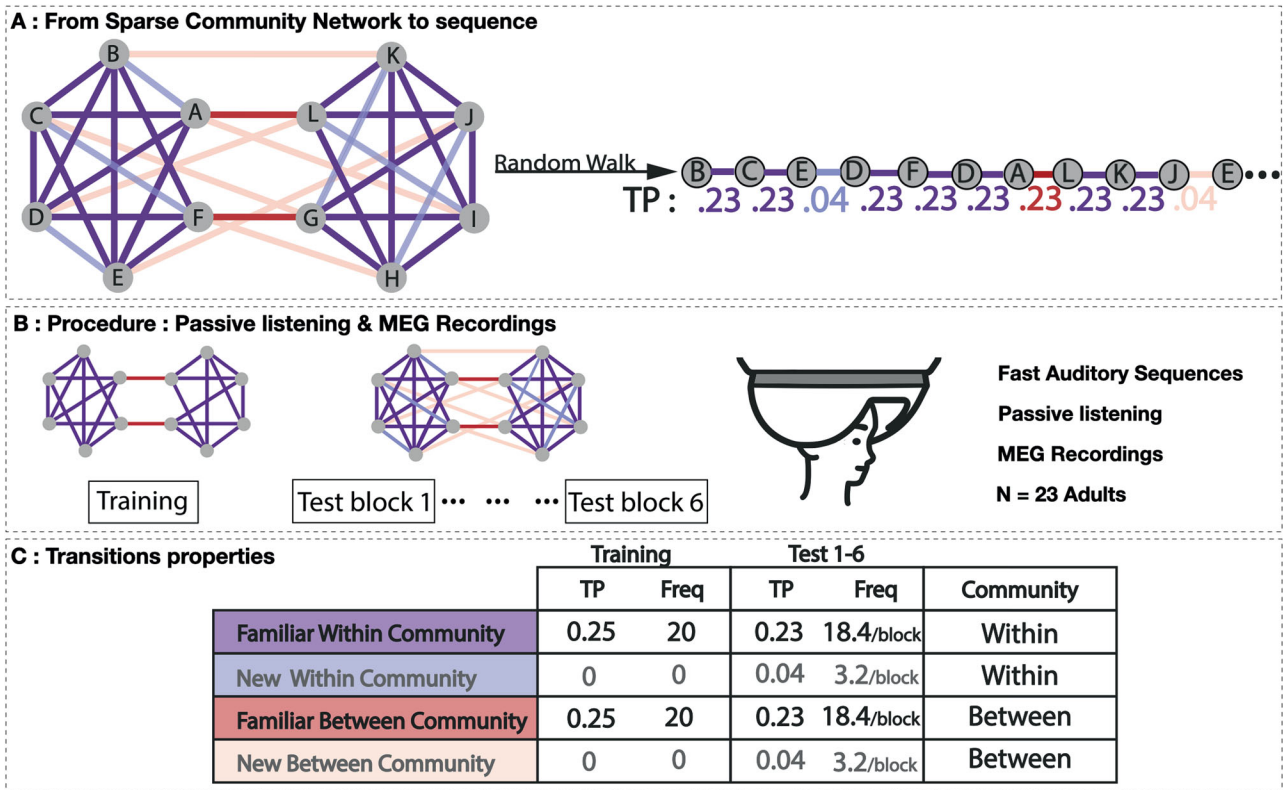
Crucially, the experiment was completely passive, and participants were unaware of the structure of the auditory sequence. They were only instructed to pay attention to the sequence of tones and to stay still while looking at a fixation cross displayed at the center of the screen to avoid noise from eye movements. The experiment lasted ~45 min, and a small break inside the MEG was possible between each block.

### Participants

Twenty-nine healthy adults came to the lab and 23 recordings (16 females; mean age, 26.58; SD = 6.1) were kept for the analyses (4 subjects were rejected due to MEG malfunction, 1 due to experimenter error during recording, and 1 scan was aborted due to subject agitation). All participants gave written informed consent prior to enrollment and received 90€ as compensation. This experiment was approved by the national ethical committee (CPP Ile-de-France III).

### MEG recordings and preprocessing

Participants performed the tasks while sitting inside an electromagnetically shielded room. The magnetic component of their brain activity was recorded with a 306-channel, whole-head MEG by Elekta Neuromag. The MEG helmet is composed of 102 triplets, each comprising one magnetometer and two orthogonal planar gradiometers. Brain signal was acquired at a sampling rate of 1,000 Hz with a hardware high-pass filter at 0.03 Hz. The data were then resampled at 250 Hz to reduce computational load. Eye movements were monitored with vertical and



**Movie 1.** Design and procedure. **A**, Example of a sparse community network for one participant. All community networks are similar in terms of properties, but New Within and New Between transitions are randomly drawn for each participant. Purple lines correspond to Familiar Within community transitions, red lines to Familiar Between community transitions, and blue and pink lines correspond, respectively, to New Within and New Between transitions. We can derive a sequence by performing a random walk into this network (click to see video of the design). Here we display an example of a test sequence derived from this structure. **B**, Experimental procedure. First, participants passively listened to a sequence from a sparse community network, in which each TP between tones was 25% (Training). Then they were presented with six 960-item test blocks obtained from the community structure graph comprising New Within and Between community transitions with low transition probabilities of 4% (light blue and pink colors on the graph). **C**, Table summarizing the local and community properties for the transitions for each condition. Each single Familiar transition is, on average, presented 18.4 times/block (20 in the Training sequence) and, therefore, has a probability of 23% to be observed (25% in the training sequence) irrespectively of staying within or switching between communities. New transitions (Within community and Between communities) have a probability of 4% in the test blocks, which implies that each single New Transition is heard 3.2 times/block on average. [View online]

horizontal EOGs and heartbeats with ECGs. Subjects' head position inside the helmet was measured for realignment at the beginning of each run with an isotrack Polhemus system from the location of four coils placed over the frontal and mastoids.

MEG signal was then preprocessed using MNE-Python pipeline with classical steps following recommendations from Jas et al. (2018) and Niso et al. (2018). We first applied Maxfilter algorithm to remove ambient noise, and signal was bandpass filtered ([0.1–30] Hz). Eye movements and heartbeats were identified and removed using PCA components' correlation with EOG and ECG measures.

To decode if a transition was within or between community, data was epoched from 100 ms before to 300 ms after tone onset. To determine how sustained was the neural representation of each tone across time, we segmented the data in 2.6 s long epochs, from 100 ms before to 2,500 ms after tone onset. Bad data, channels, and epochs were detected and removed with autoreject toolbox (Jas et al., 2017).

*Within versus Between decoding analysis*

To examine whether the brain encoded the community structure, we trained a logistic regression decoder to predict whether the transition that just occurred stayed Within a community (Familiar Within and New Within) or switched Between communities (Familiar Between and New Between). The decoder was trained on the short epochs ([−0.1, 0.3]) slightly smoothed using a sliding window (±20 ms) to enhance the signal-to-noise ratio. We used threefold cross-validation process: the decoder was trained on two-thirds of the data and tested on the remaining third of the trials. The procedure was repeated three times, corresponding to the three cross-validation folds. Each transition

had the same frequency, but Within transitions were more numerous than Between transitions, resulting in a larger total number of epochs for Within condition. We thus used the area under the ROC curve as a metric of success (ROC AUC) since it is not sensitive to such imbalance. This analysis was conducted for each time point of the epochs (Fig. 1). We also computed the decoding performance when the decoder was trained at time *t* and tested at time *t'*, to reveal the generalization across time (GAT) of the decoder, and thus the stability of the mental representation (Fig. 1). By design, the diagonal of the GAT matrix corresponds to the previously described time-by-time decoding performance.

To assess robustness, we replicated the decoding accuracy with a different metric, and we performed a decoding analysis on the whole epoch at the subject level for training and testing. This decoder simultaneously used all time points across all recording channels, providing a single accuracy value for the entire epoch. Unlike time-by-time decoding, this approach can exploit the temporal dynamic of the signal to differentiate conditions.

For the previous analysis, we pulled together the data from Familiar and New transitions. In a further analysis, we investigated whether the success of decoding the community remained possible when analyzing Familiar and New transitions separately. Therefore, we replicated the previous decoding analysis but limited it to Familiar transitions only, which had identical high local transition probabilities of 23% (Familiar Within vs Familiar Between) or to New transitions only (New Within vs New Between), which had a low TP of 4%. Note, however, that in this last case, the number of epochs was small, resulting in a low signal-to-noise ratio.

**Statistical analysis.** Statistical significance in the GAT matrix was assessed using a temporal–temporal cluster-based permutation (MNE-Python; Gramfort et al., 2013) for times between 0 and 300 ms. For the time-by-time decoder, we performed a temporal cluster permutation test in [0, 300] ms time window. Note that these two statistical tests are not independent as the time-by-time decoding corresponds to the diagonal of the GAT matrix. The whole epoch decoding gives a single decoding value per subject, we thus performed a one-way *t* test across subjects to test whether the decoding performance was significantly above chance.

#### Long-horizon associative learning estimation and linear regression

To assess the duration of the representation of a sequence item in the brain signals, we used epochs containing 10 tones (2.5 s). We trained a 12-class decoder (for the 12 tones) with balanced accuracy to decode the identity of the first tone of the epoch throughout the whole epoch. To ensure that we were decoding the sustained activity related to the first tone and not a subsequent repetition of the same tone, we removed from the analysis all epochs in which the first tone was repeated during the test window (~65% of the epochs were removed; Fig. 3A). We averaged the above chance decoding performance over the time windows (250ms), which corresponded to the interval between two consecutive items, to estimate the amount of superposition of the representations of the different elements of the sequence. We then estimated the long-horizon associative learning strength of the association for each pair ( $\hat{A}$ ), which corresponds to the sum of the TP matrix between the tones at all orders ( $A^i$ ), weighted by the overlap between item representations (Fig. 3B).

We later used the associated novelty index, defined as the negative log of this association strength, as a regressor for the MEG signal during the short epochs corresponding to the different transition types (Fig. 3D). We performed spatiotemporal cluster analysis on the  $\beta$  value associated with this linear regression to extract electrodes and times where this long-horizon associative learning estimation might significantly explain the difference in activity across conditions. We also computed the average association strength of each type of transition (Fig. 3C).

## Results

Our experiment aimed to identify the neural correlates of community structure encoding and evaluate if this learning stems from a low-level associative process or corresponds to a late and explicit discovery (Ren et al., 2022; Stiso et al., 2022). To assess the encoding of the community structure, we first decoded Within versus Between transition type. Afterwards, we characterized the temporal dynamics of the representation of each tone in the sequence in order to assess the possibility of overlapping representations that might allow long-distance associations. Based on this measured overlap, we could estimate the long-horizon associative familiarity for each transition. Finally, to determine whether this long-horizon associative learning model was indeed a plausible hypothesis, we ran a linear regression between the predicted familiarity and our data.

### Decoding Within versus Between community transitions

We first tested whether participants' mental model of the sequence encoded the community structure despite uniform transition probabilities. We thus trained and tested decoders on all tone epochs ending in a Within transition versus all tone epochs ending in a Between transition on all pairs (Familiar and New). We obtained a significant cluster ( $p < 0.05$ ) in the GAT matrix accuracy. Temporal cluster analysis on the time-by-time decoding accuracy revealed a significant cluster between 90 and 250 ms ( $p < 0.001$ ), peaking at 160 ms. Finally, the epoch-based decoding was significantly above chance ( $p < 0.01$ ; Fig. 1, Within vs Between).

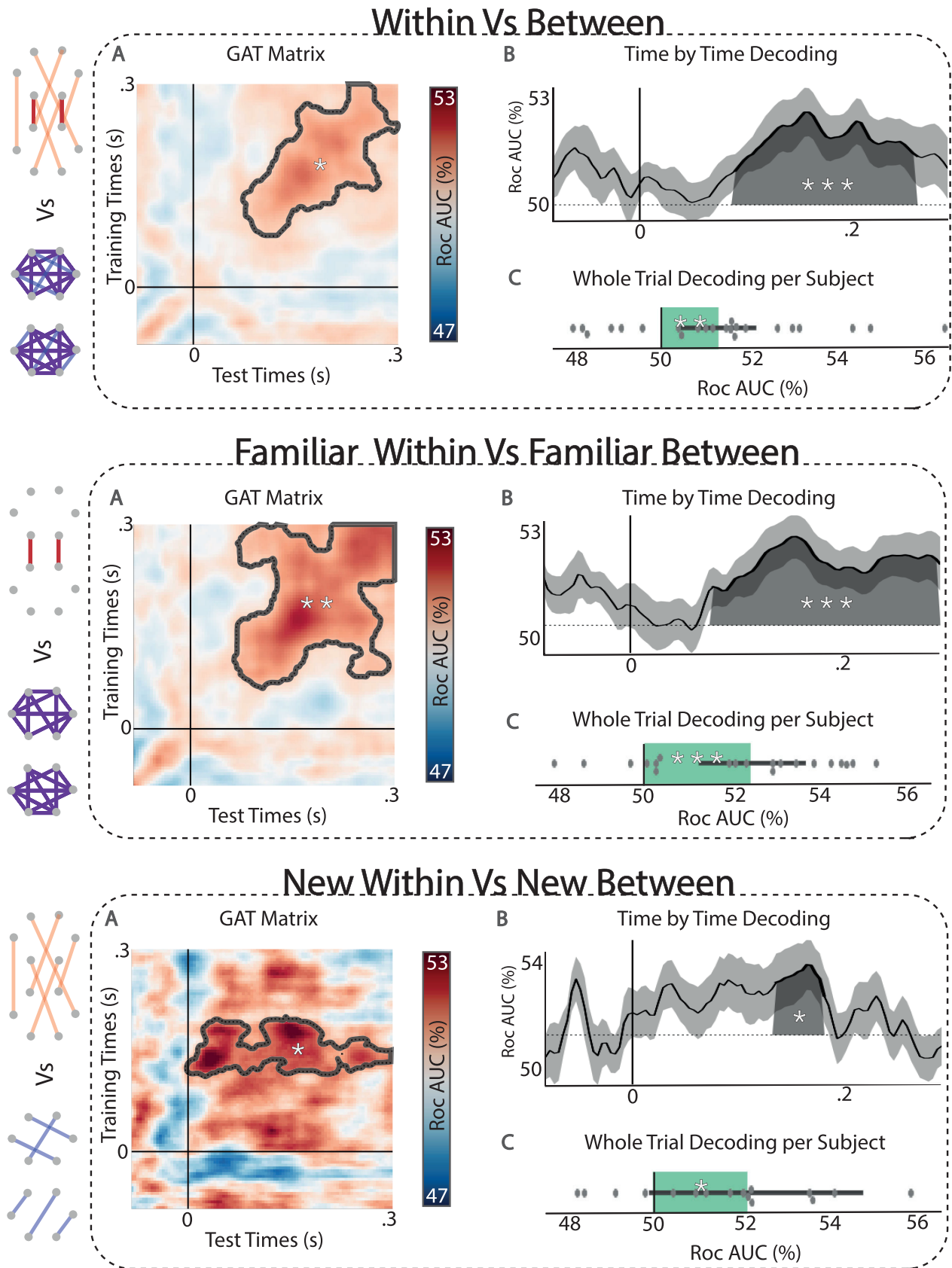
We then restricted this analysis to the Familiar transitions (Familiar Within vs Familiar Between, which corresponds to

92% of the epochs). Since Familiar Within and Familiar Between transitions have the same transition probabilities (0.23), a significant difference would then be due to a higher-order representation of the community structure. Here again, a significant cluster ( $p < 0.01$ ) was found in the GAT matrix. A temporal cluster between 80 and 280 ms was found in the time-by-time decoding ( $p < 0.001$ ) with a peak at 150 ms. Epoch-based decoding was also significantly above chance ( $p < 0.001$ ).

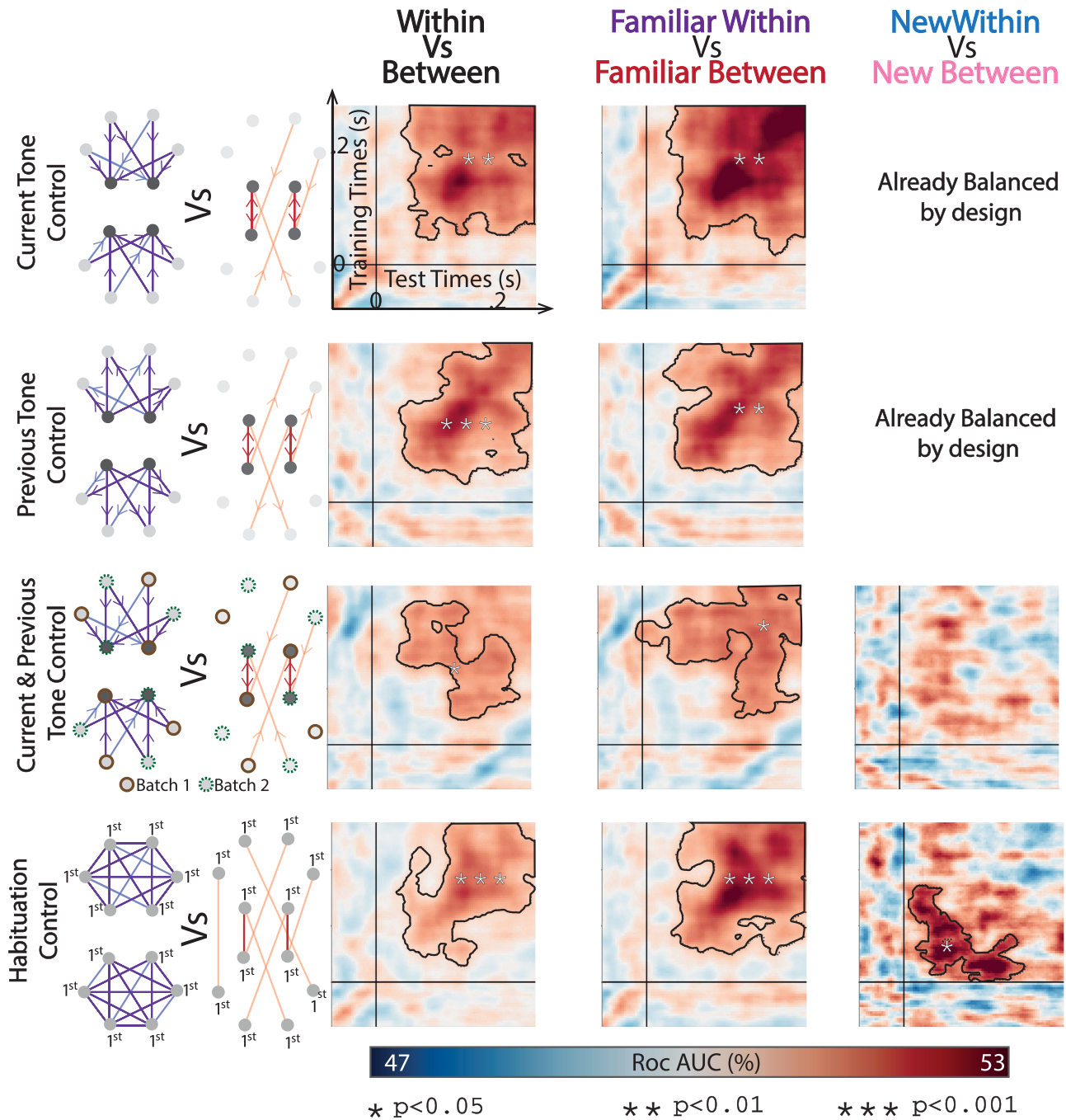
Symmetrically, we restricted the analysis to New transitions only (New Within vs New Between, which corresponds to 8% of the epochs). By design, both New Within and New Between transitions had transition probabilities of 4% so learning only local transition probabilities would predict equal unfamiliarity with both types of transition. In line with the previous results, we found a significant temporal–temporal cluster in the generalization matrix ( $p < 0.05$ ), and a significant temporal cluster in the time-by-time decoding ( $p < 0.05$ ; significant time, [130, 170] ms; peaking at 160 ms). Epoch-based decoding was also significant ( $p < 0.05$ ). Due to the much smaller number of epochs, the results were noisier.

We also computed the ERF on the gradiometers for the Familiar Within versus Familiar Between and New Within versus New Between contrasts on the [100–200] ms time window to confirm the presence of the effect found with the decoding approach. The outcomes were qualitatively comparable: a significant effect ~150 ms for the Familiar Within versus Familiar Between contrast (cluster-based permutations  $p < 0.001$ ), and a trend effect for the New Within versus New Between contrast (cluster-based permutations  $p = 0.075$ ). In both cases, the topography of the difference was compatible with an auditory response.

We performed a series of control analyses to eliminate putative low-level confounds, such as decoding success based on the identity of the current tone, the previous tone, or the pair of tones. To control for tone identity decoding, we ran the decoding analysis but restricted it to one of the four nodes at the border of a community (i.e., connected to a node of the other community, darker nodes in Fig. 2). Depending on the previous tone, these epochs could be either Familiar Within, Familiar Between, New Within, or New Between. Thus, decoding within versus between community transitions on those epochs cannot be driven by the tone identity. The same was done for epochs where the transition began with one of these four nodes (i.e., epochs where the previous node of the sequence was one of the nodes at the border of communities) to control for decoding the identity of the previous tone. We also controlled for the pair forming the transition (previous and current tone identity simultaneously): in a similar manner to the current tone control, we restricted the analysis to nodes at the borders of communities and also cross-validated the decoding on the previous tone identity. To do so, we trained and tested our decoder on different previous nodes (training on three previous nodes per community and testing on the three others, see batches in Fig. 2). This strategy was also used for the Familiar Within versus Familiar Between GAT matrix decoder. By experimental design, New Within versus New Between decoders were already balanced for current and previous tones (each node is attached to one transition of each type). Thus, we only controlled for the pair by using the cross-validation of the previous node with the same batches described above. Overall, the control analyses qualitatively and quantitatively confirmed previous results. Only the New Within versus New Between control for pair (i.e., controlling



**Figure 1.** Within versus Between community decoders on the MEG signal. Top panel, Decoders with all Within community epochs (Familiar and New) versus all Between communities epochs (Familiar and New) transitions. **A**, GAT matrix with significant cluster delineated in black. **B**, Time-by-time decoding. The shaded area indicates a significant temporal cluster. **C**, Individual performances based on whole epoch decoding: Mean decoding accuracy across subjects (green bar, one dot per subject, the black line represents standard error). Those three analyses have been replicated with Familiar only transitions (middle panel) and New only transitions (bottom panel). Community structure was encoded in each case despite the flat local TP. Stars represent significance of the statistical tests (\* $p < 0.05$ , \*\* $p < 0.01$ , \*\*\* $p < 0.001$ ).



**Figure 2.** Control analyses for the results presented in Figure 2. For each decoder, we controlled for the current tone, the previous tone, and the pair (both current and previous tone simultaneously). We also controlled for habituation due to temporal proximity between tones. All the analyses qualitatively and quantitatively confirmed previous results except the New Within versus New Between control analysis that did not reach significance, probably because of the small number of epochs.

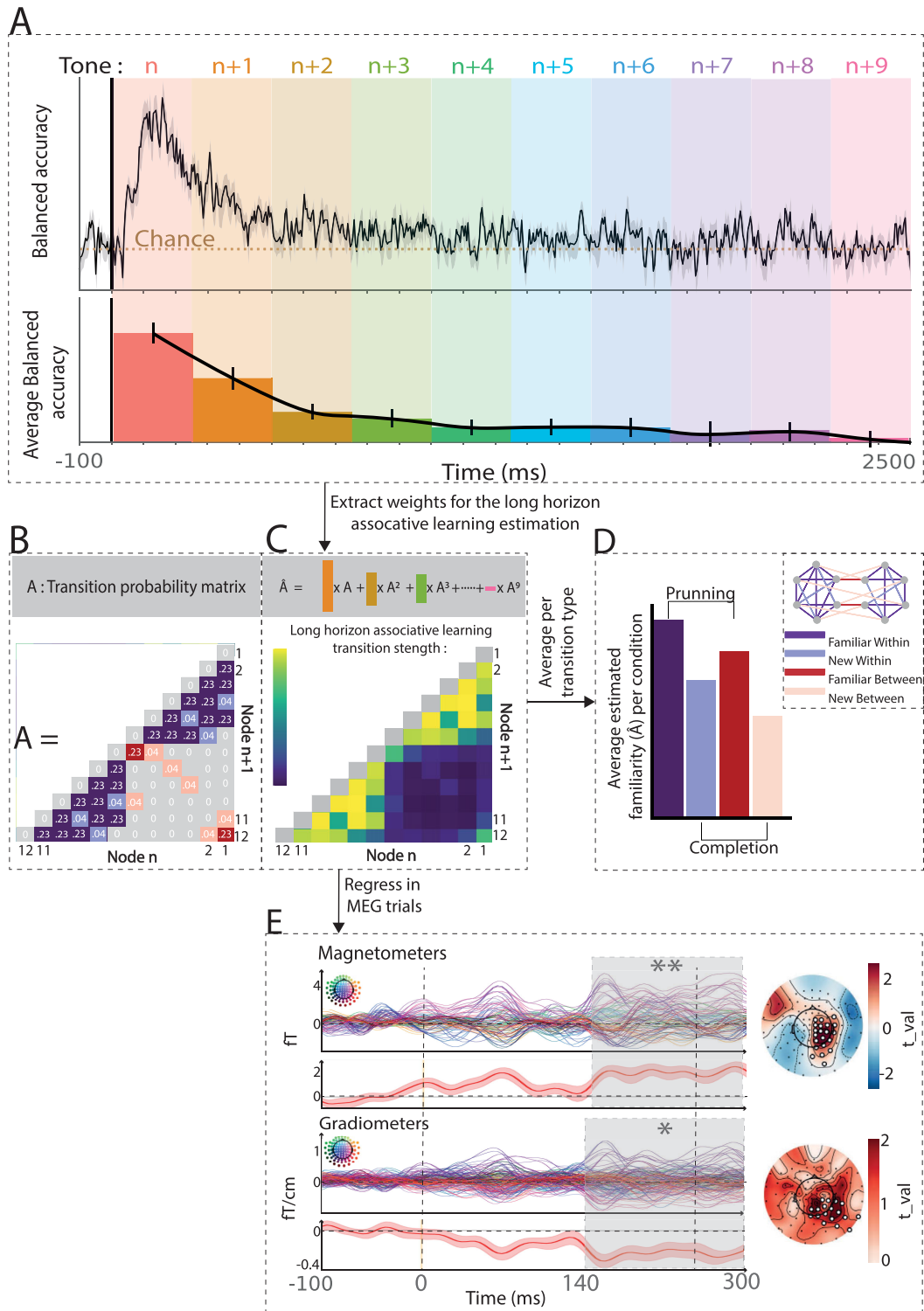
both previous and current tone identity simultaneously) analysis did not reach significance, probably due to the small number of epochs in this analysis (only 8% of the data was used in this last control).

Within versus Between community transitions decoding could also rely on a habituation effect. Indeed, if the sequence remains within a community, a particular sound might be repeated multiple times within a short span, causing habituation. However, if the sequence shifts from one community to another, the same sound is less likely to be repeated in a short time, thus preventing habituation. Therefore, this differential habituation effect could drive the Within versus Between decoder. To rule

out this alternative hypothesis, we restricted the analysis to the first appearance of each tone after a community change. Thus, close repetitions of tones of the same community are avoided in the data used for this decoder. Despite a decrease in the number of epochs, the decoding accuracy of those controls was still significant for all conditions. All generalization matrices are shown in Figure 2.

#### Long-horizon associative learning estimation

We tested here the hypothesis that long-horizon associative learning (associative learning over several consecutive and non-consecutive items) can support the encoding of network



**Figure 3.** Associative learning estimation and fit on MEG data. **A**, Top, Decoding performance of the first item of the sequence across time (2.5 s window). Shaded colors indicate the stimulus onset asynchrony (SOA) between each tone of the sequence. The dotted line shows the chance level. Bottom, Decoding performance averaged over the duration of each tone and the following ISI. Error bars present the standard error across subjects. It takes  $\sim 8$  items for the decoder of the first tone to converge to chance level. **B**, Matrix of exact transition probabilities ( $A$ ) associated with the graph underlying the sequence. Familiar transitions are associated with 23% transition probabilities and New with 4% transition probabilities (Movie 1). Impossible transitions have a null TP. **C**, Estimation of the long-horizon associative learning strength for each transition. Based on the decoder (panel **A**), we estimated the overlap between nonadjacent elements of the sequence (average decoder accuracy during SOA of item  $n+i$ ). We then computed the associative learning strength ( $\hat{A}$  matrix) for each pair of elements as the sum of the different transitional probability orders ( $A^i$ ), weighted by the overlap between item representations. **D**, Average of the long-horizon associative learning strength per condition. Pruning (Familiar Within > Familiar Between) and completion (New Within > New Between) effects are consistent with behavioral results (Benjamin et al., 2023a) and with the decoding performance obtained in Fig 1. **E**, Regression coefficient for the estimated long-horizon associative novelty ( $-\log(\hat{A})$ ) for each MEG sensor. Significant time windows are shown in shaded areas and significant sensors are indicated on the  $t$ -map topographies by the white dots. These were obtained with a spatiotemporal cluster-based permutation analysis. The red line below the sensors value represents the time course of the average regression value on the sensors of the significant cluster.

structure. This concept builds on Hebb's principle of strengthening the link between co-occurring events. Nonetheless, instead of focusing solely on learning adjacent pairs, we proposed a broader approach that allows connections to be established over longer distances. In our experiment, this long-horizon associative learning implies that the mental representation of each tone is sustained for a sufficient duration to allow several tones to overlap (Endress, 2010) and thus enable associations through more successive tones. According to this model, it is predicted that the representation of each tone should decrease following an exponential profile. To test this hypothesis, we quantified the overlap between the representations of item  $n$  and item  $n + i$ . In fact, this provides a good estimator of the weight of the nonadjacent TP of order  $i$ .

To estimate the overlap between brain representations of different items of the sequence, we determined how long the representation of each item was seen in brain activity. To do so, we split the data into 10-item-long sequences (i.e., 2.5 s) with no repetition of the first tone in the sequence. We train a 12-class decoder on each time point to predict the identity of the first tone. Decoding performance is shown in Figure 3. We averaged the above chance decoding performance over the time windows corresponding to the interval between two consecutive items. We observed an exponential-like decrease in performance that reached 0 after ~8 sequence items (Fig. 3A). It shows that the overlap enabling associative learning might thus include long-horizon dependencies of up to eight items.

We estimated the long-horizon associative learning strength of each pair of tones. To do so, we computed the sum of the different transitional probability orders weighted by the overlap between item representations as estimated from the decoding performances (Fig. 3B). This gave us a  $12 \times 12$  symmetrical matrix of learning familiarity for each pair (Fig. 3C). Finally, we averaged this measure of Familiarity for each condition type (Fig. 3C) and obtained a result that is consistent with the pruning effect (difference between Familiar Within vs Familiar Between transitions) and the completion effect (difference between New Within and New Between transitions) as discussed in Benjamin et al. (2023a).

### Long-horizon associative learning accounts for epoch variability

To test the neural predictions of long-horizon associative learning, we correlated brain signals with the estimated associative learning strength of each transition (Fig. 3D). We performed a linear regression between the brain signal after each tone and the novelty effect produced by each transition. Unlike most studies of sequence learning, where the novelty is calculated solely from local transition probabilities, we computed it here as the negative log of the long-horizon associative learning strength. This calculation takes into account several orders of adjacent and nonadjacent transition probabilities whose weights have been computed based on the overlap of brain representations estimated by our tone decoder (Fig. 3A–D). A spatiotemporal cluster permutation test revealed a significant cluster (Fig. 3E) in the magnetometers (right centro-occipital; time, [150; 290] ms;  $p$  value < 0.01) that was replicated in the gradiometers (right centro-occipital; time, [140; 300] ms;  $p$  value < 0.05).

Furthermore, the observed clusters were still significant when the negative log of the adjacent transition probabilities was introduced as a supplementary regressor ( $ps < 0.05$  for both magnetometers and gradiometers clusters). However, while the strong correlation between the TP matrix and the long-horizon

associative model makes it hard to directly disentangle those two models solely based on this regression analysis, it does complement the decoding analyses nicely.

## Discussion

In this study, our aim was to determine whether local statistical learning and structure learning in sequences are governed by the same cognitive process or by distinct processes. Learning local statistics is often described as an associative process, while network learning is usually seen as an abstract map representation. Previous studies exploring network learning have used explicit paradigms, revealing late brain signatures consistent with top-down or frontal activity (Ren et al., 2022; Stiso et al., 2022). However, based on a modeling approach, we proposed in our previous behavioral study that low-level associative learning strategies might support both local and high-order statistical scales (Benjamin et al., 2023a). Thus, this hypothesis predicts that learning sequence structure does not require an explicit representation and may instead rely on automatic and rapid (~150 ms) mismatch responses, similar to those observed after the violation of local transition probabilities.

### Network learning results from a low-level bottom-up computation

To test these predictions, we presented participants with a passive learning task using rapid auditory sequences. We showed that the structure properties of the sequence were rapidly decodable from the participants' brain recordings (~[100–250] ms after tone onset). The timing of this response, as early as 150 ms after the information became available, aligns with the rapid deviant responses (MMN in EEG) observed in learning based on violation of transitional probabilities (Todorovic and de Lange, 2012; Maheu et al., 2019). Since the transition probabilities between tones were uniform and the walk within the network was random, prediction could not be based on high-level top-down expectation. This early and automatic response (150 ms after the transition) challenges the notion of abstract and explicit calculations as prerequisites for learning such structures. In addition, our analyses revealed a similar effect when the decoding analysis was restricted to new transitions (New Within vs New Between) and to familiar transitions (Familiar Within vs Familiar Between), suggesting an automatic generalization of the community structure beyond sensory evidence. This result provides a neural underpinning for the behavioral observations we previously reported, indicating that participants accurately assess the familiarity of transitions based on their congruence with network structure, even when these transitions were not encountered during training.

### Long-horizon associative learning as a plausible implementation for FEMM

In our previous study, we hypothesized that the FEMM could effectively explain adult behavioral performance. This model aggregates the different orders of statistical regularities (adjacent and nonadjacent) into a single quantity. In this study, we showed that this model can be readily implemented through a simple associative learning mechanism relying on Hebb's principle (Hebb, 1949; Benjamin et al., 2023a). In the context of structure learning, this principle would imply a sustained mental representation of each tone for a sufficient duration to enable the overlapping of several elements despite the temporal distance. We thus predicted the representation of each tone to exhibit an



exponential decay profile. A rapid decay of tone information would limit associations to short distances, while a slower decay would facilitate the formation of long-horizon dependencies and, therefore, the extraction of the underlying structure. Thus, this exponential decay acts as a balance between local relevance and generalization.

To test this idea, we estimated the duration of the representation of each tone, performing a decoding analysis of tone identity. The identity of a tone was decodable during the presentation of the subsequent eight tones, with a decoding performance exponentially decaying over subsequent tones. This profile provided an estimation of the number of elements simultaneously represented at a given time. Consequently, it allowed us to quantitatively assess the strength of each tone pair in the heard sequence (Fig. 3C). We found that these weights accurately accounted for the results of the Within versus Between decoders, encompassing both Familiar and New transitions (Fig. 3D). Moreover, this estimated strength significantly correlated with neural activity, aligning with the timing of the automatic deviant response (Todorovic and de Lange, 2012; Maheu et al., 2019). This result provides compelling evidence for the rapid encoding of structure through bottom-up processes compatible with associative learning strategies.

However, it is worth noting that an alternative implementation of the same metric is theoretically possible. Simple pairwise association learning, in combination with a transitivity property, would also predict similar learning. In fact, if participants solely learn pairs (e.g., A–C and C–D), transitivity of this learning can strengthen the A–D pair, even if not explicitly presented. Considering this transitivity with similar exponentially decreasing weights would be mathematically equivalent to our model while not strictly requiring a sequential presentation of the structure. Although, we cannot definitively rule out this alternative implementation of the same metric, our findings suggest that sequential presentation is crucial to have an overlap between successive items representations, enabling Hebbian associative learning. It is also important to acknowledge that associative learning might not be the sole mechanism contributing to network structure learning, particularly in cases where explicit detection is required from participants. Abstract representations of hippocampal maps (Constantinescu et al., 2016) or frontal maps (Stiso et al., 2022) might also play a role in such tasks (Schapiro et al., 2016, 2017; Garvert et al., 2017). Intracranial recordings conducted during local statistical learning paradigms have revealed that multiple brain regions, including cortical areas and hippocampus, can simultaneously represent the same structure while carrying different information (Henin et al., 2021).

### Difference between implicit passive listening and explicit structure learning

Thus, converging results provide evidence that associative learning supports the perception of the community structure in the present experiment. Long-horizon associative learning strength significantly accounted for the variance in brain signals (Fig. 3E). Moreover, the pruning and completion effects found with decoders (Fig. 1) can easily be explained by the same mechanism (Fig. 3D). However, it is worth noting that the results from our previous behavioral study do not entirely align with the current ones. Specifically, in the present experiment, the representation of tones exhibited a more rapid decrease (exponential decrease factor 0.52) as compared with its estimation in our previous behavioral study (factor 0.058, ~10 times lower). This

discrepancy suggests that participants in the current experiment might be less inclined to generalize the underlying structure.

Several factors might explain this difference. Firstly, the generalization factor estimation in the MEG experiment may be noisier due to the small number of participants (23 vs several hundred in the behavioral study). Since the trade-off between generalization and accuracy may vary among individuals (Lynn et al., 2020), on the one side, group level estimation with 23 subjects is limited, and on the other side at the individual level, it is difficult to measure this trade-off due to the data variability. A larger sample size with multiple sessions per subject would be necessary to obtain a reliable estimation of the generalization factor at the individual level. Secondly, it is possible that associative learning represents the implicit component of this task (Andringa and Rebuschat, 2015), followed subsequently by an explicit decision-making process involving higher level prefrontal regions. This second step might facilitate the abstraction of the structure by labeling each community as distinct (Koechlin et al., 2003; Koechlin and Jubault, 2006). This dual process could explain why explicit behavioral tasks (Lynn et al., 2020; Benjamin et al., 2023a) exhibit a better generalization factor compared with our implicit MEG task. The same explanation may account for the late signatures of top-down activity reported by Ren et al. (2022) who used a slow and explicit task. To further explore this hypothesis, a direct comparison of passive and active learning of such networks while monitoring the representations in the auditory cortex, the hippocampus, and the lateral prefrontal cortex would be necessary.

### Conclusion

The aim of the present study was to uncover the neural mechanism underlying network learning. We proposed the sparse community paradigm as a way of combining local statistical learning and network learning in a single sequence. Previous behavioral studies have shown that a mathematical model (FEMM) accurately captures human learning. Here, we add that the behavioral pattern described by the FEMM is compatible with certain associative learning principles. Indeed, thanks to time-by-time decoding of the brain state associated with a tone, we observed an exponential decay in the tone representation across eight elements. Using this estimate of mental representations' dynamics, we estimated the strength of each network transition. This estimate significantly correlated with our data. The present study provides novel insights into the mechanism underlying network learning and highlights the importance of brain dynamics in the understanding of sequence learning. Further investigations in different experimental conditions (explicit vs implicit), over different tone and ISI durations, with different populations (non-human primates), and during early development are necessary to better characterize this learning ability.

### References

- Andringa S, Rebuschat P (2015) New directions in the study of implicit and explicit learning: an introduction. *Stud Second Lang Acquis* 37:185–196.
- Benjamin L, Dehaene-Lambertz G, Fló A (2021) Remarks on the analysis of steady-state responses: spurious artifacts introduced by overlapping epochs. *Cortex* 142: 379–388.
- Benjamin L, Fló A, Al Roumi F, Dehaene-Lambertz G (2023a) Humans parsimoniously represent auditory sequences by pruning and completing the underlying network structure. *eLife* 12:e86430.
- Benjamin L, Fló A, Palu M, Naik S, Melloni L, Dehaene-Lambertz G (2023b) Tracking transitional probabilities and segmenting auditory sequences are dissociable processes in adults and neonates. *Dev Sci* 26:e13300.

- Benjamin L, Zang D, Flo A, Qi Z, Su P, Zhou W, Wang L, Wu X, Gui P, Dehaene-Lambertz G (2024) The role of conscious attention in statistical learning: evidence from patients with impaired consciousness.
- Boros M, Magyari L, Török D, Bozsik A, Deme A, Andics A (2021) Neural processes underlying statistical learning for speech segmentation in dogs. *Curr Biol* 31:5512–5521.e5.
- Constantinescu AO, Jill O, Behrens TEJ (2016) Organizing conceptual knowledge in humans with a gridlike code. *Science* 352:1464–1468.
- Dehaene S, Al Roumi F, Lakretz Y, Planton S, Sablé-Meyer M (2022) Symbols and mental programs: a hypothesis about human singularity. *Trends Cogn Sci* 26:751–766.
- Dehaene S, Meyniel F, Wacongne C, Wang L, Pallier C (2015) The neural representation of sequences: from transition probabilities to algebraic patterns and linguistic trees. *Neuron* 88:2–19.
- Endress AD (2010) Learning melodies from non-adjacent tones. *Acta Psychol* 135:182–190.
- Endress AD, Johnson SP (2021) When forgetting fosters learning: a neural network model for statistical learning. *Cognition* 213:104621.
- Fló A, Benjamin L, Palu M, Dehaene-Lambertz G (2022) Sleeping neonates track transitional probabilities in speech but only retain the first syllable of words. *Sci Rep* 12:4391.
- Garvert MM, Dolan RJ, Behrens TEJ (2017) A map of abstract relational knowledge in the human hippocampal–entorhinal cortex. *eLife* 6:1–20.
- Gramfort A, et al. (2013) MEG and EEG data analysis with MNE-Python. *Front Neurosci* 7:267.
- Hebb DO (1949) *The organization of behavior: a neuropsychological theory*. Mahwah, NJ: L. Erlbaum Associates.
- Henin S, Turk-Browne NB, Friedman D, Liu A, Dugan P, Flinker A, Doyle W, Devinsky O, Melloni L (2021) Learning hierarchical sequence representations across human cortex and hippocampus. *Sci Adv* 7:1–13.
- James LS, Sun H, Wada K, Sakata JT (2020) Statistical learning for vocal sequence acquisition in a songbird. *Sci Rep* 10:1–18.
- Jas M, Engemann DA, Bekhti Y, Raimondo F, Gramfort A (2017) Autoreject: automated artifact rejection for MEG and EEG data. *NeuroImage* 159:417–429.
- Jas M, Larson E, Engemann DA, Leppäkangas J, Taulu S, Hämäläinen M, Gramfort A (2018) A reproducible MEG/EEG group study with the MNE software: recommendations, quality assessments, and good practices. *Front Neurosci* 12:530.
- Karuz EA, Kahn AE, Bassett DS (2019) Human sensitivity to community structure is robust to topological variation. *Complexity* 2019:1076.
- Koechlin E, Jubault T (2006) Broca's area and the hierarchical organization of human behavior. *Neuron* 50:963–974.
- Koechlin E, Ody C, Kouneiher F (2003) The architecture of cognitive control in the human prefrontal cortex. *Science* 302:1181–1185.
- Lynn CW, Kahn AE, Nyema N, Bassett DS (2020) Abstract representations of events arise from mental errors in learning and memory. *Nat Commun* 11:2313.
- Maheu M, Dehaene S, Meyniel F (2019) Brain signatures of a multiscale process of sequence learning in humans. *eLife* 8:1–24.
- Mark S, Moran R, Parr T, Kennerley SW, Behrens TEJ (2020) Transferring structural knowledge across cognitive maps in humans and models. *Nat Commun* 11:1–12.
- Marr D (1982) *Vision: a computational investigation into the human representation and processing of visual information*. Cambridge, MA: MIT Press.
- Niso G, et al. (2018) MEG-BIDS, the brain imaging data structure extended to magnetoencephalography. *Sci Data* 5:180110.
- Ren X, Zhang H, Luo H (2022) Dynamic emergence of relational structure network in human brains. *Prog Neurobiol* 219:102373.
- Saffran JR, Aslin RN, Newport EL (1996) Statistical learning by 8-month-old infants. *Science* 274:1926–1928.
- Schapiro AC, Rogers TT, Cordova NI, Turk-Browne N, Botvinick MM (2013) Neural representations of events arise from temporal community structure. *Nat Neurosci* 16:486–492.
- Schapiro AC, Turk-Browne NB, Botvinick MM, Norman KA (2017) Complementary learning systems within the hippocampus: a neural network modelling approach to reconciling episodic memory with statistical learning. *Philos Trans R Soc B Biol Sci* 372:372.
- Schapiro AC, Turk-Browne NB, Norman KA, Botvinick MM (2016) Statistical learning of temporal community structure in the hippocampus. *Hippocampus* 26:3–8.
- Stiso J, et al. (2022) Neurophysiological evidence for cognitive Map formation during sequence learning. *eNeuro* 9:ENEURO.0361-21.2022.
- Todorovic A, de Lange FP (2012) Repetition suppression and expectation suppression are dissociable in time in early auditory evoked fields. *J Neurosci* 32:13389–13395.
- Toro JM, Trobalón JB (2005) Statistical computations over a speech stream in a rodent. *Percept Psychophys* 67:867–875.