

1 **Multiscale integration organizes hierarchical computation in human**
2 **auditory cortex**

3

4 *Sam V Norman-Haignere^{1,2}, Laura K. Long^{1,3}, Orrin Devinsky^{4,5}, Werner Doyle^{5,6}, Ifeoma
5 Irobunda⁷, Edward M. Merricks⁷, Neil A. Feldstein⁷, Guy M. McKhann⁷, Catherine A. Schevon⁷,
6 Adeen Flinker^{4,5}, Nima Mesgarani^{1,3,8}

7

8 1. Zuckerman Mind, Brain, Behavior Institute, Columbia University

9 2. HHMI Postdoctoral Fellow of the Life Sciences Research Foundation

10 3. Doctoral Program in Neurobiology and Behavior, Columbia University

11 4. Department of Neurology, NYU Langone Medical Center

12 5. Comprehensive Epilepsy Center, NYU Langone Medical Center

13 6. Department of Neurosurgery, NYU Langone Medical Center

14 7. Department of Neurology, Columbia University Medical Center

15 8. Department of Electrical Engineering, Columbia University

16 *Corresponding and lead author

17

18 Keywords: integration period, auditory cortex, timescale, intracranial EEG,
19 electrocorticography, natural sounds

20 **Abstract**

21

22 To derive meaning from sound, the brain must integrate information across tens (e.g.
23 phonemes) to hundreds (e.g. words) of milliseconds, but the neural computations that enable
24 multiscale integration remain unclear. Prior evidence suggests that human auditory cortex
25 analyzes sound using both generic acoustic features (e.g. spectrotemporal modulation) and
26 category-specific computations, but how these putatively distinct computations integrate
27 temporal information is unknown. To answer this question, we developed a novel method to
28 estimate neural integration periods and applied the method to intracranial recordings from
29 human epilepsy patients. We show that integration periods increase three-fold as one ascends
30 the auditory cortical hierarchy. Moreover, we find that electrodes with short integration periods
31 (~50-150 ms) respond selectively to spectrotemporal modulations, while electrodes with long
32 integration periods (~200-300 ms) show prominent selectivity for sound categories such as
33 speech and music. These findings reveal how multiscale temporal analysis organizes
34 hierarchical computation in human auditory cortex.

35 Time is the fundamental dimension of sound, and temporal integration is thus fundamental to
36 audition. To recognize a complex structure like a word, the brain must integrate information
37 across a wide range of timescales from tens (e.g. phonemes) to hundreds (e.g. syllables) of
38 milliseconds (**Fig S1**)¹. But how human auditory cortex accomplishes this feat is unclear.

39

40 One prominent hypothesis posits that short and long-term temporal structure are analyzed
41 asymmetrically across the two hemispheres, with the left hemisphere integrating over short
42 timescales, and the right hemisphere integrating over long timescales²⁻⁴. Another influential
43 hypothesis is that the auditory cortex integrates across time hierarchically, with short-term
44 structure analyzed bilaterally in primary auditory cortex and longer-term structure analyzed in
45 non-primary regions⁵⁻⁷. This question remains unresolved, despite intensive debate over two
46 decades, because the integration period of human cortical regions is unknown.

47

48 Understanding temporal integration is critical for understanding how important sound
49 categories like speech and music are processed^{2,6,8}. While prior studies have revealed non-
50 primary neural populations selective for speech and music⁹⁻¹³, little is known about how these
51 neural populations integrate information in speech and music. One possibility is that category-
52 selective neural populations integrate over many timescales in order to code category-specific
53 structure at short^{14,15} (e.g. phonemes) and long⁸ timescales (e.g. syllables and words; **Fig S1**).
54 Alternatively, short-term structure might be analyzed by general-purpose acoustic
55 representations in primary auditory cortex¹⁶ and then integrated over long timescales to form
56 category-specific neural representations in non-primary regions.

57

58 Here, we test these hypotheses by developing a novel method for measuring neural integration
59 periods. Integration periods are often defined as the time window when stimuli alter the neural
60 response^{17,18}. Although this definition is simple and general, there is no simple and general
61 method to estimate integration periods. Many methods exist for inferring linear integration
62 periods with respect to a spectrogram^{15,19-21}, but human cortical responses exhibit prominent
63 nonlinearities particularly in non-primary regions. Flexible, nonlinear models are challenging to
64 fit given limited neural data^{20,22}, and even if one succeeds, it is not obvious how to measure
65 the model's integration period. Methods for assessing temporal modulation selectivity^{6,23,24} are
66 insufficient, since a neuron could respond to fast modulations over a long window or to a
67 complex structure like a word that is poorly described by its modulation content. Finally,
68 temporal scrambling can reveal selectivity for naturalistic temporal structure^{12,18,25}, but many
69 regions in auditory cortex show no difference between intact and scrambled sounds.

70

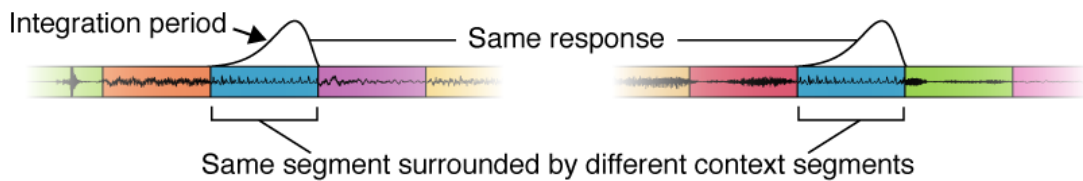
71 To overcome these limitations, we developed a method that directly estimates the time window
72 when stimuli alter a neural response (the temporal context invariance or TCI paradigm; **Fig 1**).
73 We present sequences of natural stimuli in a random order such that the same segment occurs
74 in different contexts. While context has many meanings²⁶, here we simply define context as
75 the stimuli which surround a segment. If the integration period is shorter than the segment
76 duration, there will be a moment when it is fully contained within each segment. As a
77 consequence, the response to each segment will be unaffected by the surrounding segments.
78 We can therefore estimate the integration period by determining the minimum segment
79 duration needed to achieve a context invariant response.

80

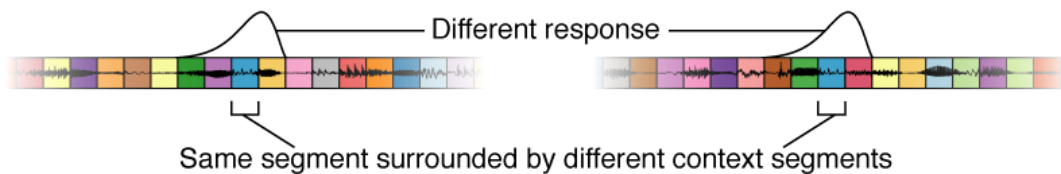
81 TCI does not make any assumptions about the type of response being measured. As a
82 consequence, the method is applicable to sensory responses from any modality, stimulus set,
83 or recording method. We applied TCI to intracranial EEG (iEEG) recordings collected from
84 patients undergoing surgery for intractable epilepsy. Such recordings provide a rare
85 opportunity to measure human brain responses with spatiotemporal precision, which is
86 essential to studying temporal integration. We used a combination of depth and surface

87 electrodes to record from both primary regions in the lateral sulcus as well as non-primary
88 regions in the superior temporal gyrus (STG), unlike many iEEG studies that have focused on
89 just the lateral sulcus²⁷ or STG^{15,19}. The precision and coverage of our recordings was
90 essential to revealing how the human auditory cortex integrates across multiple timescales.
91

Segment duration > Integration period



Segment duration < Integration period



92
93 **Fig 1. Temporal context invariance (TCI) paradigm.** Schematic of the paradigm used to
94 measure integration periods. Segments of natural stimuli are presented using two different
95 random orderings. As a consequence, the same segment occurs in two different contexts
96 (different surrounding segments). If the segment duration is longer than the integration period
97 (top panel), there will be a moment when the integration period is fully contained within each
98 segment. As a consequence, the response at that moment will be unaffected by the surrounding
99 context segments. If the segment duration is shorter than the integration period (bottom panel),
100 the integration period will always overlap the surrounding context segments, and they can
101 therefore alter the response. The goal of the TCI paradigm is to estimate the minimum segment
102 duration needed to achieve a context invariant response. This figure plots waveforms for an
103 example sequence of segments that share the same central segment. Segment boundaries are
104 demarcated by colored boxes. The hypothesized integration period is plotted above each
105 sequence at the moment when it best overlaps the shared segment.
106

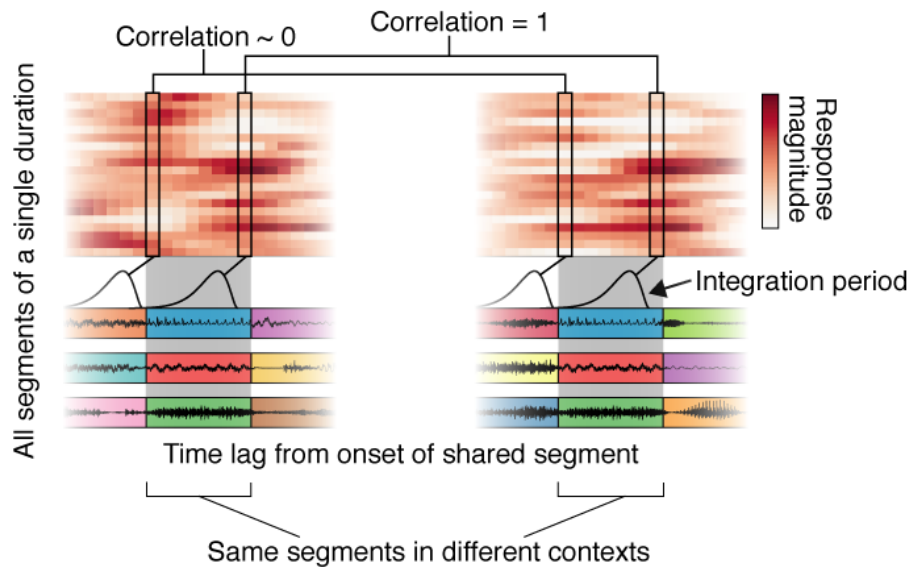
107 Results

108
109 We recorded iEEG responses to sequences of natural sound segments that varied in duration
110 (from 31 ms to 2 sec in octave steps). For each segment duration, we created two 20-second
111 sequences, each with a different random ordering of the same segments (concatenated using
112 cross-fading to avoid boundary artifacts). Segments were excerpted from 10 natural sounds
113 (**Table S1**), selected to be diverse so they differentially drive responses throughout auditory
114 cortex. The same natural sounds were used for all segment durations, which limited the
115 number of sounds we could test given the limited time with each patient; but our key results
116 were robust across the sounds tested (see *Anatomical organization* for the results of all
117 robustness analyses). Because our goal was to characterize integration periods during natural
118 listening, we did not give subjects a formal task. To encourage subjects to listen to the sounds,
119 we asked them to occasionally rate how scrambled the last stimulus sequence was (shorter
120 segment durations sound more scrambled; if patients were in pain or confused we simply
121 asked them to listen).
122

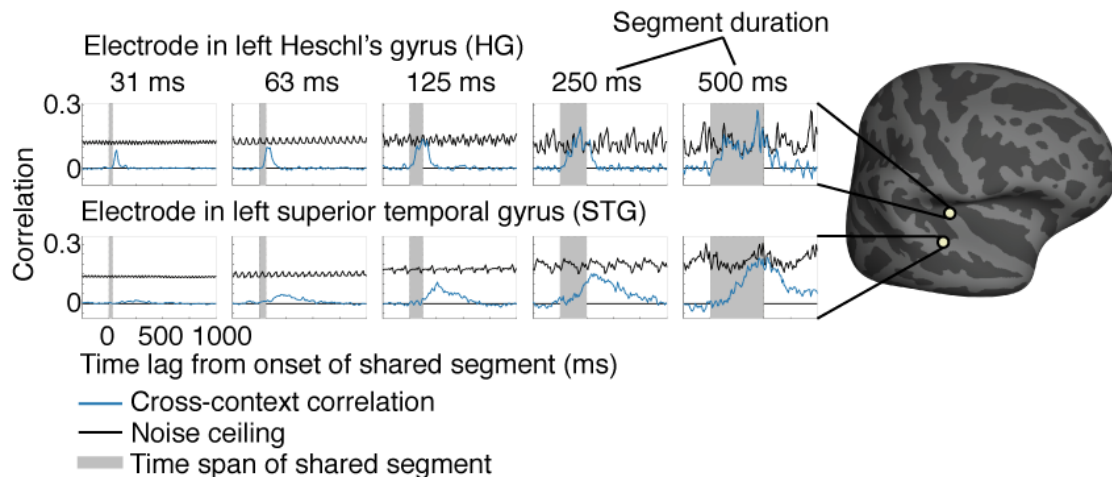
123 **Assessing context invariance via the cross-context correlation.** We measured the
124 broadband gamma power of each electrode to each sound sequence, which is thought to
125 approximately reflect aggregate neural activity in a local region^{28,29} (70-140 Hz; results were
126 robust to the frequency range used). For each electrode, we aligned its response to all

127 segments of a given duration in a matrix, which we refer to as the segment-aligned response
 128 (SIR) matrix (**Fig 2a**). Each row of the SIR matrix contained the response timecourse
 129 surrounding a single segment, aligned to segment onset. Different rows thus correspond to
 130 different segments and different columns correspond to different lags relative to segment
 131 onset. We computed two versions of the SIR matrix using the two different contexts for each
 132 segment, extracted from the two different sequences. The central segment was the same
 133 between the contexts, but the surrounding segments were different.
 134

a Schematic of cross-context correlation



b Cross-context correlation for example electrodes



135 **Fig 2. Cross-context correlation.** **a**, Schematic of the analysis used to assess context invariance
 136 for a single electrode and segment duration. See text for description. **b**, The cross-context
 137 correlation (blue line) and noise ceiling (black line) are shown for two example electrodes from
 138 the left hemisphere of the same patient, one in Heschl's gyrus (HG, top panel) and one in the
 139 superior temporal gyrus (STG, bottom panel). Each plot shows a different segment duration. The
 140 gray region shows the time interval when the shared segment was present (i.e. the gray region in
 141 panel a). The STG electrode required longer segment durations for the cross-context correlation
 142 to reach the noise ceiling, and the build-up of the cross-context correlation with lag was slower
 143 for the STG electrode.
 144

145
 146 Our goal was to assess if there was a lag when the response was the same across contexts.
 147 We instantiated this idea by correlating corresponding columns across SIR matrices from
 148 different contexts (the "cross-context correlation", schematized in **Fig 2a**). At segment onset

149 (lag=0), the cross-context correlation should be near zero, since the integration period must
150 overlap the preceding segments, which were random across contexts. As time progresses, the
151 integration period should start to overlap the shared segment, and the cross-context
152 correlation should increase. Critically, if the integration period is less than the segment
153 duration, there should be a lag where the integration period is fully contained within the shared
154 segment, and the response should thus be the same, yielding a correlation of 1 modulo noise.
155 To correct for noise, we measured the test-retest correlation when the context was the same,
156 which provides a noise ceiling for the cross-context correlation.

157
158 The shorter segments tested in our study were created by subdividing the longer segments.
159 As a consequence, we could also consider cases where a segment was a subset of a longer
160 segment and thus surrounded by its natural context, in addition to the case described so far
161 when a segment is surrounded by random other segments. Since our analysis requires that
162 the two contexts being compared are different, one of the two contexts must always be
163 random, but the other context can be random or natural. In practice, we found similar results
164 using random and natural contexts, and thus pooled across both types of context for maximal
165 statistical power (see *Anatomical organization* for results comparing random and natural
166 contexts).

167
168 We plot the cross-context and noise ceiling for segments of increasing duration for two
169 example electrodes from the same subject: an electrode in left posteromedial Heschl's gyrus
170 (HG) and one in the left superior temporal gyrus (STG) (**Fig 2b**). The periodic variation in the
171 noise ceiling is an inevitable consequence of correlating across a fixed set of segments (see
172 *Cross-context correlation* in the Methods for an explanation). For the HG electrode, the cross-
173 context correlation started at zero and quickly rose. Critically, for segment durations greater
174 than or equal to 125 milliseconds, there was a lag where the cross-context correlation equaled
175 the noise ceiling, indicating a context invariant response. For longer segments (250 or 500
176 ms), the cross-context correlation remained yoked to the noise ceiling for an extended duration
177 indicating that the integration period remained within the shared segment for an extended time
178 period. This pattern is what one would expect for an integration period that is ~125
179 milliseconds, since stimuli falling outside of this window have little effect on the response.

180
181 By comparison, the results for the STG electrode suggest a much longer integration period.
182 Only for segment durations of approximately 500 milliseconds did the cross-context correlation
183 approach the noise ceiling, and its build-up and fall-off with lag was considerably slower. This
184 pattern is what one would expect for a longer integration period, since it takes more time for
185 the integration period to fully enter and exit the shared segment. Virtually all electrodes with a
186 reliable response to sound exhibited a similar pattern, but the segment duration and lag
187 needed to achieve an invariant response varied substantially (**Fig S2** shows 20 representative
188 electrodes). This observation indicates that auditory cortical responses have a meaningful
189 integration period, outside of which responses are largely invariant, but the extent of this
190 integration period varies across the cortex.

191
192 **Model-estimated integration periods.** In theory, one could estimate the integration period
193 extent as the shortest segment duration for which the peak of the cross-context correlation
194 exceeds some fraction of the noise ceiling. This approach, however, would be noise-prone
195 since a single noisy data point at one lag and segment duration could alter the estimated
196 integration period. To overcome this issue, we used a model to infer the integration period that
197 best-predicted the cross-context correlation across all lags and segment durations.

198

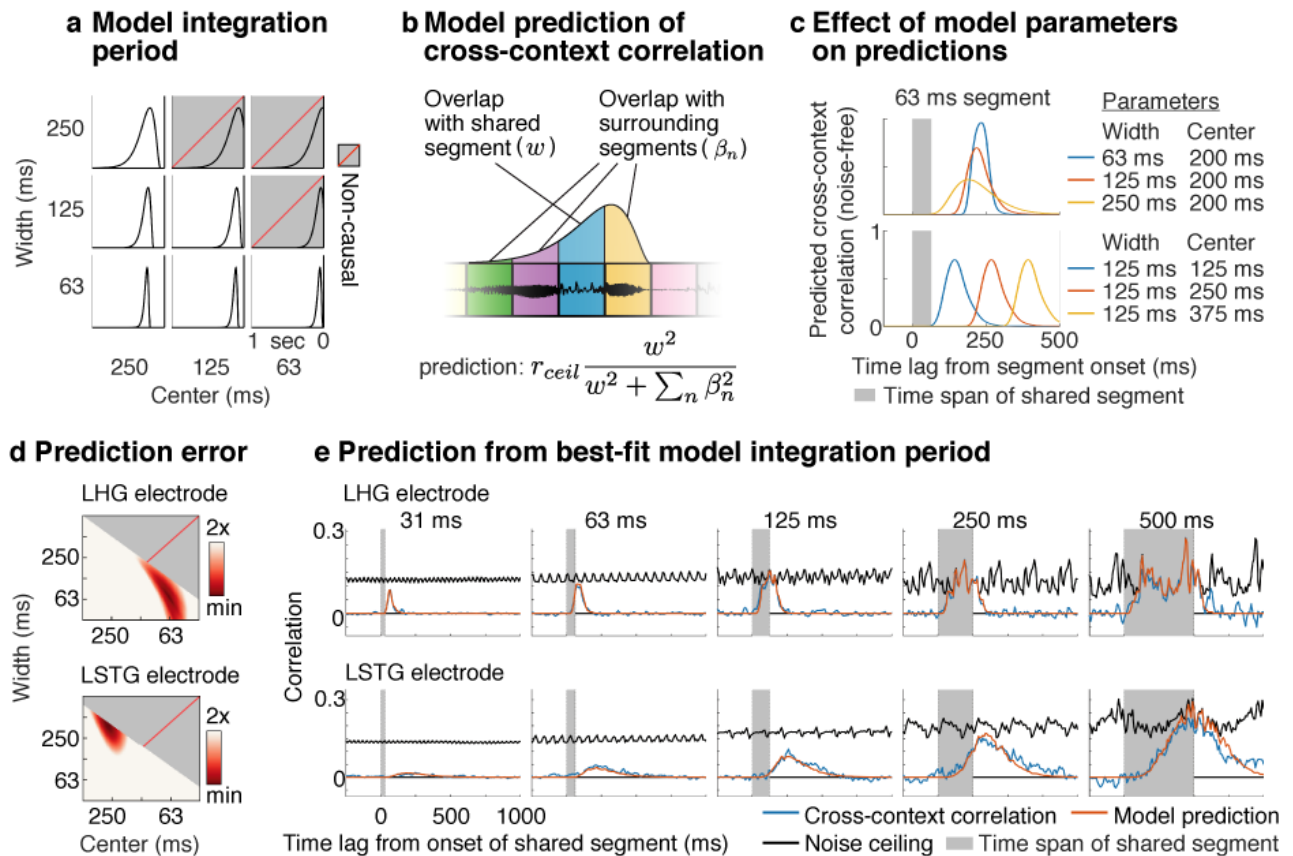


Fig 3. Model-estimated integration periods. **a**, Temporal integration periods were modeled using a Gamma-distributed window. The width and center of the model integration period were varied, excluding combinations of widths and centers that resulted in a non-causal window (gray boxes with dashed red line). **b**, Schematic showing how the cross-context correlation was predicted from the model integration period. For each segment duration and lag, we measured how much the integration period overlapped the shared central segment (w , blue segment) vs. all surrounding context segments (β_n , yellow, purple, and green segments). The cross-context correlation should reflect the fraction of the response variance due to the shared segment, multiplied by the noise ceiling (r_{ceil}). The variance due to each segment is given by the squared overlap with the model integration period. **c**, Illustration of how the integration period width (top panel) and center (bottom panel) alter the model's prediction for a single segment duration (63 milliseconds). Increasing the width lowers and stretches-out the predicted cross-context correlation, while increasing the center shifts the cross-context correlation to later lags. **d**, The prediction error for model windows of varying widths and centers for the example electrodes from **Figure 2b**. Redder colors indicate lower error. **e**, The measured and predicted cross-context correlation for the best-fit integration period with lowest error (same format as **Fig 2b**).

We modeled temporal integration periods using a Gamma-distributed window, which is a standard, unimodal distribution commonly used to model integration periods (**Fig 3a**)³⁰. We varied the width and center of the model integration period, excluding combinations of widths and centers that resulted in a non-causal window since this would imply the response depends upon future stimuli. The width of the integration period is the key parameter we would like to estimate, and was defined as the smallest interval that contained 75% of the window's mass. The center of the integration period was defined as the window's median and reflects the overall delay between the integration period and the response. We also varied the window shape from more exponential to more bell-shaped, but found the shape had little influence on the results (see *Anatomical organization*).

The cross-context correlation depends on the degree to which the integration period overlaps the shared segment vs. the surrounding context segments (**Fig 2a**). We therefore predicted

230 the cross-context correlation by measuring the overlap between the model integration period
231 and each segment, separately for all lags and segment durations (**Fig 3b**). The equation used
232 to predict the cross-context correlation from these overlap measures is shown in **Figure 3b**
233 and described in the legend. A formal derivation is given in the Methods.

234

235 **Figure 3c** illustrates how changing the width and center of the model integration period alters
236 the predicted correlation. Increasing the width lowers the peak of the cross-context correlation,
237 since a smaller fraction of the integration period overlaps the shared segment at the moment
238 of maximum overlap. The build-up and fall-off with lag is also more gradual for wider integration
239 periods since it takes longer for the integration period to enter and exit the shared segment.
240 Increasing the center simply shifts the cross-context correlation to later lags, since the delay
241 is longer, but the width is unchanged.

242

243 We varied the model parameters and calculated the error between the measured and predicted
244 cross-context correlation (**Fig 3d,e**). For the example HG electrode, the cross-context
245 correlation was best-predicted by an integration period with a narrow width (68 ms) and early
246 center (64 ms) compared with the STG electrode, which was best-predicted by a wider and
247 more delayed integration period (375 ms width, 273 ms center). These results validate our
248 qualitative observations and provide us with a quantitative estimate of each electrode's
249 integration period.

250

251 **Anatomical organization.** We identified 190 electrodes with a reliable response to sound
252 across 18 patients (test-retest correlation > 0.1 ; $p < 10^{-5}$ via a permutation test across sound
253 sequences; 128 left hemisphere; 62 right hemisphere). From these electrodes, we created a
254 map of integration widths and centers, discarding a small fraction of electrodes (5%) where the
255 model predictions were not highly significant ($p < 10^{-5}$ via a phase-scrambling analysis) (**Fig**
256 **4a**). This map was created by localizing each electrode on the cortical surface, and aligning
257 each subject's brain to a common anatomical template. By necessity, we focus on group
258 analyses due to the sparse, clinically-driven coverage in any given patient. Most sound-
259 responsive electrodes were in and around the lateral sulcus and STG, as expected^{11,15}.

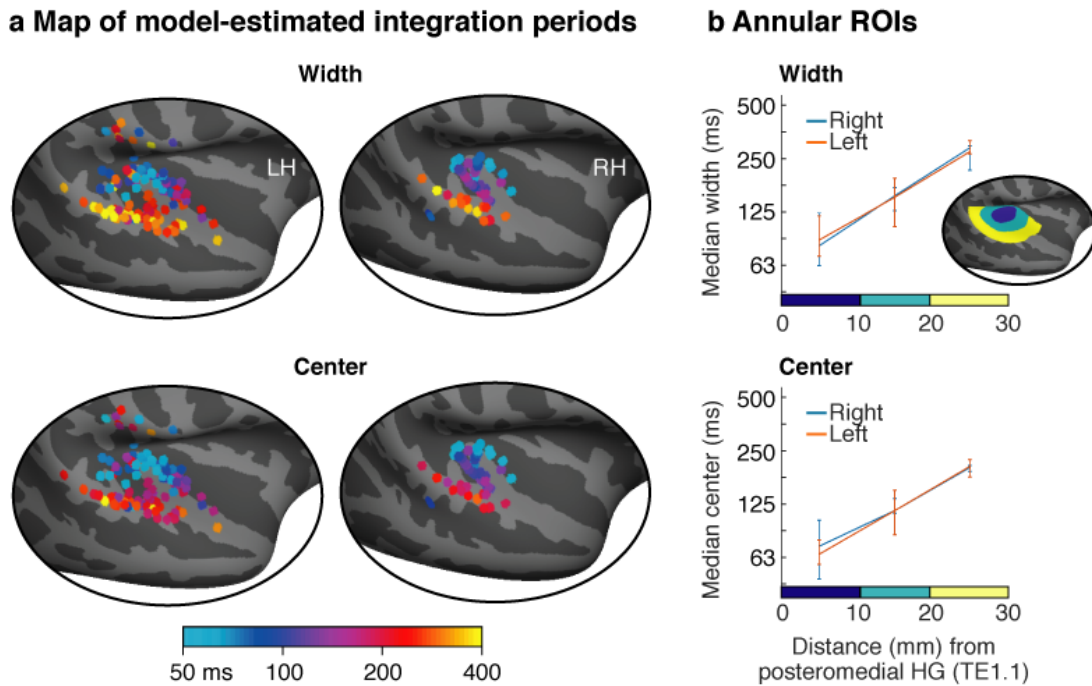
260

261 These maps revealed a clear anatomical gradient: integration widths and centers increased
262 substantially from primary regions near posteromedial HG to non-primary regions near STG.
263 We quantified this trend by binning electrodes into anatomical regions-of-interest (ROIs) based
264 on their distance to posteromedial HG (**Fig 4b**)³¹. This analysis revealed a three-fold increase
265 in integration widths and centers from primary to non-primary regions (median integration
266 width: 84, 152, 281 ms; median integration center: 68, 115, 203 ms; $p < 0.001$ via a
267 bootstrapping analysis across subjects comparing the nearest and farthest bins). By contrast,
268 there was no difference in integration widths or centers between the two hemispheres either
269 when averaging across all ROIs or comparing individual ROIs (all $ps > 0.74$). These findings
270 were robust across the specific sounds tested (**Fig S3**), the type of context used to assess
271 invariance (random vs. natural; **Fig S4**), the shape of the model window (**Fig S5**), and the
272 frequency range used to measure broadband gamma (**Fig S6**). These results demonstrate
273 human auditory cortex integrates across time hierarchically, with substantially wider and more
274 delayed integration periods in higher-order regions, but no difference between hemispheres.

275

276 Across all electrodes, we found that integration centers scaled approximately linearly with
277 integration widths (**Fig S7**). In part as a consequence of this observation, we found that
278 integration centers were relatively close to the minimum possible for a causal window even
279 when not explicitly constrained to be causal (**Fig S7**) (integration centers were on average
280 46% greater than the minimum for a Gamma-distributed window). Since the integration center

281 reflects the delay between the integration period and the response, this finding suggests that
282 auditory cortex responds to sounds about as quickly as possible given the integration period³².
283

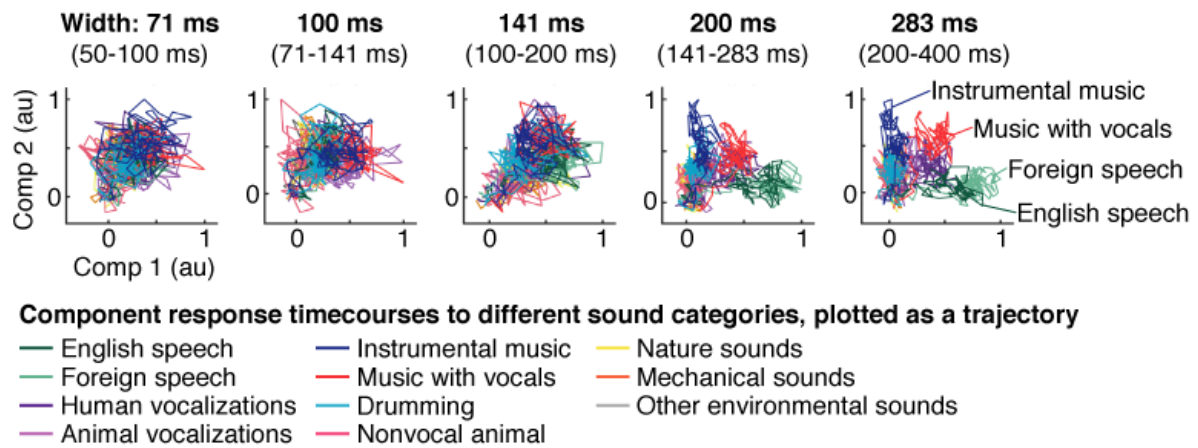


284 **Fig 4. Anatomy of model-estimated integration periods.** **a**, Map of integration widths (top)
285 and centers (bottom) for all electrodes with a reliable response to sound. **b**, Electrodes were
286 binned into ROIs based on their distance to a common anatomical landmark of primary auditory
287 cortex (posteromedial Heschl's gyrus, TE1.1). This figure plots the median integration width and
288 center across the electrodes in each bin. Inset shows the ROIs for one hemisphere. Error bars
289 plot one standard error of the bootstrapped sampling distribution across subjects.
290
291

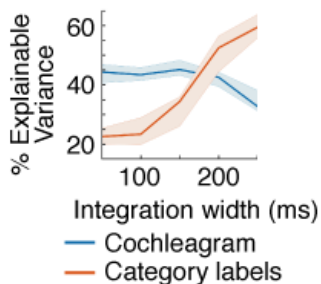
292 **Category selective responses are limited to electrodes with long integration periods.**
293 What is the consequence of hierarchical temporal integration for the analysis of important
294 sound categories like speech and music? While prior studies have revealed non-primary
295 neural populations that respond selectively to sound categories like speech and music^{9–13}, it
296 is unclear how these neural populations integrate information in speech and music. A priori it
297 seemed possible that speech and music-selective responses might have diverse integration
298 periods since sound categories like speech and music have unique structure at many
299 timescales. For example, the median duration of speech phonemes in the popular TIMIT
300 corpus is 64 milliseconds while syllables and words are typically hundreds of milliseconds (the
301 median duration of multiphone syllables and multisyllable words is 197 and 479 milliseconds,
302 respectively) (**Fig S1**). But the hierarchy revealed by our integration period maps suggests an
303 alternative hypothesis: that category-selective responses are limited to neural responses with
304 wide integration periods. We sought to directly test this hypothesis, and if true, determine the
305 shortest integration periods at which category-selective responses are present.
306

307 To assess category selectivity, we ran a separate experiment in a subset of 11 patients, where
308 we measured responses to a larger set of 119 natural sounds, drawn from 11 categories (listed
309 in **Fig 5a**). We grouped the electrodes from these patients based on their integration width in
310 octave intervals (shown in **Fig 5a**), pooling across both hemispheres because we had fewer
311 electrodes (104) and because integration periods (**Fig 4**) and category-selective responses<sup>10–
312 12</sup> are similar across hemispheres. We then used several different analyses to measure the
313 degree of category selectivity in each electrode group.
314

a Category-selective response components at different integration periods



b Prediction accuracy



c Speech & music selectivity

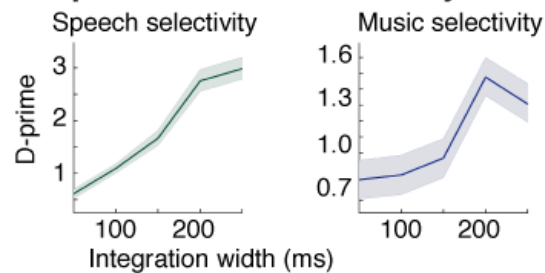


Fig 5. Category selectivity at different integration period widths. Responses were measured in a subset of patients to a larger collection of 119 natural sounds from 11 different sound categories, listed in panel a. Electrodes from these patients were grouped based on the width of their integration period in octave intervals. We used several different analyses to assess the degree of category-selectivity in each group. **a**, The responses from each group were projected onto the two components that exhibited the greatest category selectivity. We plot the timecourses for these two components as a 2D trajectory after averaging across the sounds from each category. Category selectivity, if present, will cause the trajectories to separate from each other. **b**, The accuracy of category labels (red line) and cochleagrams (blue line) in predicting electrode responses as a function of the integration width. This figure plots the squared correlation (noise-corrected) between the measured and predicted response for each feature set. **c**, The selectivity of each electrode group for either speech (English and foreign) or music (instrumental, vocal and drumming), measured along the component that exhibited the greatest speech or music selectivity. Selectivity was measured as the separation (noise-corrected d-prime) between responses to sounds from the target category (speech or music) compared with all other sounds. Independent sounds were used to estimate components and measure their response. Error bars plot one standard error of the bootstrapped sampling distribution.

First, we used component methods to visualize any category selectivity present in the electrodes¹¹. Specifically, we projected the responses from each group onto the two components that showed the greatest category selectivity (**Fig 5a**) (estimated using cross-validated linear discriminant analysis). We plot the timecourse of these two components as a 2D trajectory after averaging across the sounds from each category. Category selectivity, if present, will cause the trajectories for different categories to separate from each other. This analysis revealed that only electrodes with wide integration periods, above ~200 milliseconds, show robust category selectivity. Similar results were obtained when analyzing the top two principal components without any optimization for category selectivity (**Fig S8**), which demonstrates that category selectivity is a prominent feature of cortical responses with wide integration periods.

346 To quantify these trends, we measured how accurately we could linearly predict the response
347 of each electrode from either a cochleagram or binary category labels (**Fig 5b**). Cochleagrams
348 are similar to spectrograms but are computed using filters designed to mimic the pseudo-
349 logarithmic frequency resolution of cochlear filtering³⁰. This analysis thus provides an estimate
350 of the fraction of the response that can be predicted using a linear spectrotemporal receptive
351 field^{15,19,20}. The category labels indicated the membership of each sound in each category for
352 all timepoints with sound energy above a minimum threshold. Prediction accuracies were
353 noise-corrected using the test-retest reliability of the electrode responses, which provides an
354 upper bound on the fraction of the response explainable by any model²⁰.

355
356 We found that the prediction accuracy of the category labels more than doubled as integration
357 widths increased, with a relatively sharp increase at ~200 milliseconds ($p < 0.001$ via
358 bootstrapping across subjects). In contrast, the prediction accuracy of cochlear features
359 decreased ($p < 0.05$), yielding a significant interaction between feature type and integration
360 width ($p < 0.001$). This finding confirms our observation that responses become substantially
361 more category-selective as integration periods widen.

362
363 We note the absolute prediction accuracies were modest for both the cochleagram and
364 category labels, never exceeding more than 45% and 60% of the explainable response
365 variance, respectively (as expected). This fact illustrates the utility of having a model-
366 independent way of estimating integration periods, since even our best-performing models fail
367 to explain a large fraction of the response, and the best-performing model can vary across
368 electrodes.

369
370 The component trajectories (**Fig 5a**) suggested that selectivity for both speech and music
371 increase as integration periods widen. To directly test this hypothesis, we separately measured
372 the degree of speech and music selectivity at different integration widths. Selectivity was
373 measured as the average separation (d -prime) between responses to speech or music vs. all
374 other sounds, along the component that showed the greatest speech or music selectivity (**Fig**
375 **5c**; **Fig S9** shows the response timecourse of these components to each category). Speech
376 sounds comprised both English and foreign speech, since we found they produced similar
377 response trajectories, consistent with prior work showing that speech-selective responses in
378 STG are not driven by linguistic meaning^{11,12}. Non-speech sounds comprised all categories
379 except vocal music, which produced an intermediate response along the speech-selective
380 component, likely due to speech in the vocals (**Fig S9**). Music sounds included instrumental
381 music, vocal music and drumming, since we have found that all three produce higher-than-
382 average responses in music-selective regions³³. We found that selectivity for both speech and
383 music increased as integration periods widen ($p < 0.001$ via bootstrapping). This increase was
384 more prominent for speech than music (note the different y-axes in **Fig 5c**), plausibly because
385 speech-selective responses are more prominent particularly in the posterior/middle STG
386 where many of our electrodes were located^{11,12,15}. But for both speech and music, there was
387 a marked increase in selectivity starting at ~200 milliseconds. This finding demonstrates that
388 increased category selectivity is a general feature of neural responses with wide integration
389 periods that applies to multiple sound categories.

390 391 **Discussion**

392
393 Our study reveals how the human auditory cortex integrates across multiple timescales when
394 processing natural sounds. Our findings resolve a longstanding debate by showing that
395 auditory cortex integrates across time hierarchically, with substantially longer integration
396 periods in non-primary regions, but no difference between hemispheres. Our results also
397 reveal the significance of hierarchical temporal integration for the analysis of important sound

398 categories like speech and music. In particular, we found that category-selective neural
399 responses are restricted to electrodes with long integration periods above ~200 milliseconds.
400 This finding suggests that short-term structure in speech and music is analyzed by general-
401 purpose acoustic representations in primary auditory cortex and then integrated over long
402 timescales to form category-specific representations in non-primary regions.

403

404 These findings were enabled by a novel method that makes it possible to estimate the
405 integration period of any sensory response. Unlike prior methods, TCI makes no assumptions
406 about the type of response being measured; it simply estimates the time window when stimuli
407 alter the neural response. As a consequence, the method should be applicable to any modality,
408 stimulus set, or recording method. We applied TCI to intracranial recordings from epilepsy
409 patients, using surface and depth electrodes placed throughout human auditory cortex. As a
410 consequence, we were able to determine how anatomically and functionally distinct regions of
411 the human brain collectively integrate across multiple timescales.

412

413 ***Implications for multiscale temporal integration.*** Temporal integration plays a central role
414 in most theories and models of auditory processing^{2-5,7}. But it has remained unclear how the
415 human auditory cortex integrates across multiple timescales, because there has been no
416 general method for estimating neural integration periods.

417

418 Hemispheric models posit that the left and right hemisphere are specialized for integrating
419 across distinct timescales^{2,5}, in part to represent the distinctive temporal structure of sound
420 categories like speech and music⁶. Some of the clearest evidence for hemispheric
421 specialization comes from recent studies that have filtered-out temporal modulations at
422 different rates in natural stimuli, and shown that removing fast temporal modulations has a
423 greater impact on responses in the left auditory cortex, particularly when processing
424 speech^{6,23}. However, the time window that a neuron integrates over cannot be determined
425 from its modulation tuning. For example, a neuron could respond selectively to fast temporal
426 modulations over a long temporal window or to a complex structure (e.g. word) that is poorly
427 described by its modulation content.

428

429 Another common proposal is that the auditory cortex integrates across time hierarchically^{3,4,7}.
430 Early evidence for hierarchical temporal organization came from the observation that “phase-
431 locking” slows from the periphery to the cortex^{34,35}, which implies that neurons encode
432 temporal modulations via changes in firing rate rather than responding at particular modulation
433 phases. But the integration period of a nonlinear response cannot be inferred from the
434 presence or absence of phase-locking. For example, a neuron could integrate over several
435 cycles of a modulation and then respond at the predicted peak or trough of the oscillation
436 (phase-locking) or modulate its firing rate based on the oscillation period (rate coding).

437

438 In auditory cortex, single-unit recordings in ferrets have revealed a slight temporal broadening
439 of linear spectrotemporal receptive fields in non-primary vs. primary auditory cortex (36 vs. 33
440 ms between PEG and A1)³⁶. But the overall integration period cannot be inferred from these
441 estimates since cortical responses exhibit prominent nonlinearities^{20,22}, particularly in non-
442 primary regions³¹ and particularly when responding to natural sounds³⁷. In humans, several
443 studies have revealed selectivity for naturalistic temporal structure in non-primary regions.
444 Examples include selectivity for phonotactic structure above and beyond tuning for individual
445 phonemes^{32,38,39}, and selectivity for natural speech compared with temporally scrambled
446 speech^{12,18,25}. But again, integration periods cannot be inferred from this type of data. For
447 example, primary regions respond similarly to intact and scrambled stimuli, even for stimuli
448 that are scrambled at timescales well below the integration period of the neural response (e.g.

449 30 milliseconds)¹². As a consequence, there has been no way to test if primary and non-
450 primary regions differ in their integration period, and if so, by how much.

451
452 Because of these methodological limitations, there are no estimates of integration periods in
453 human auditory cortex, making it impossible to test how auditory cortex integrates over multiple
454 timescales. Our study thus resolves a longstanding debate by showing that multiscale
455 temporal integration is predominantly hierarchical and not hemispheric. These results do not
456 imply that there is no hemispheric organization for modulation tuning, since as already noted,
457 one cannot infer the integration period of a neural response from its modulation tuning, or vice
458 versa. Indeed, integration periods are useful specifically because they abstract away from the
459 particular features which drive a neural response. As a consequence, integration periods can
460 be used to compare any two brain regions, even if they respond to very different stimulus
461 features, like primary and non-primary auditory cortex.

462
463 The hierarchical organization of temporal integration periods appears analogous to the
464 hierarchical organization of spatial receptive fields in visual cortex^{40,41}, which suggests there
465 might be general principles that underlie this type of organization. For example, both auditory
466 and visual recognition become increasingly challenging at large temporal and spatial scales,
467 because the dimensionality of the input grows exponentially. Hierarchical multiscale analysis
468 may help overcome this exponential expansion by allowing sensory systems to gradually
469 recognize large-scale structures as combinations of smaller-scale structures (e.g. a face from
470 nose and eyes, or a word from several phonemes) rather than attempting to recognize large-
471 scale structures directly from the high-dimensional input^{3,4,7}.

472
473 **Implications for the analysis of sound categories.** Sound categories like speech and music
474 have unique acoustic structure at many temporal scales^{1,3,8,42}, from tens to hundreds of
475 milliseconds (**Fig S1**). Prior studies have revealed non-primary neural populations selective
476 for important sound categories like speech and music^{9–13}. But very little is known about how
477 information in speech and music is integrated in these neural populations. A prior fMRI study
478 used a scrambling technique called “quilting” to show that speech-selective regions respond
479 selectively to intact temporal structure up to about 500 milliseconds in duration¹². But this study
480 was only able to identify a single analysis timescale across all of auditory cortex, likely because
481 scrambling is a coarse manipulation and fMRI a coarse measure of the neural response.

482
483 We were able to identify a broad range of integration periods from tens to hundreds of
484 milliseconds, and could thus test whether category-selective responses were present at both
485 short and long integration periods. Our results indicate that category-selective responses are
486 only robustly present at integration periods above ~200 milliseconds, which corresponds to
487 about the duration of a multi-phone syllable (**Fig S1**) (the duration of musical structures is less
488 stereotyped and thus harder to assess). This finding suggests that short-term structures (e.g.
489 phonemes) are analyzed by general-purpose acoustic representations in primary auditory
490 cortex (e.g. spectrotemporal receptive fields)¹⁶, and then integrated over long timescales to
491 form category-specific representations in non-primary regions. This finding does not imply that
492 speech-selective regions are insensitive to short-term structure such as phonemes, but rather
493 that speech-selective responses respond to larger-scale patterns, such as phoneme
494 sequences, consistent with recent work on phonotactics^{32,38,39}.

495
496 We found that both speech and music selectivity increased for integration periods greater than
497 ~200 milliseconds (**Fig 5c**), which is perhaps surprising given that speech and music have
498 distinctive temporal structure^{2,43}. This result might be explained by the fact that speech and
499 music-selective responses emerge at a similar point in the cortical hierarchy^{10,11}, just beyond
500 primary auditory cortex. If integration periods are predominantly organized hierarchically, as

501 suggested by our data, regions with a similar hierarchical position might exhibit similar
502 integration periods even if they respond to different stimuli.

503
504 **Limitations and extensions.** As with any method, our results could depend upon the stimuli
505 tested. We tested a diverse set of natural sounds with goal of characterizing responses
506 throughout auditory cortex using ecologically relevant stimuli. Because time is short when
507 working with surgical patients, we could only able to test a small number of sounds, but found
508 that our key findings were robust to the sounds tested (**Fig S3**). Nonetheless, it will be
509 important in future work to test whether and how integration periods change for different
510 stimulus classes.

511
512 Another key question is whether temporal integration periods reflect a fixed property of the
513 cortical hierarchy or whether they are shaped by attention and behavioral demands. In our
514 study, we did not give subjects a formal task because our goal was to measure integration
515 periods during natural listening without any particular goal or attentional focus. Future work
516 could explore how behavioral demands shape temporal integration periods by measuring
517 integration periods in the presence or absence of focused attention to a short-duration (e.g.
518 phoneme) or long-duration (e.g. word) target⁴⁴.

519
520 Temporal integration periods indicate the time window in the stimulus that a neural response
521 is sensitive to, which is distinct from the time window in the neural response that best encodes
522 a property of the stimulus, sometimes referred to as the “encoding window”^{17,45–47}. For
523 example, the encoding window for a slow temporal modulation could be quite long, even if the
524 neural integration period is quite short. Encoding windows thus reflect a complex mixture of
525 the neural integration period and the temporal properties of the stimulus being encoded.

526
527 A given neural response might effectively have multiple integration periods. For example,
528 neural responses are known to adapt their response to repeated sounds on the timescale of
529 seconds⁴⁸ to minutes⁴⁹ and even hours⁵⁰, suggesting a form of long-term memory⁵¹. TCI
530 measures the integration period of responses that are reliable across repetitions, and as a
531 consequence, TCI will be insensitive to response characteristics that change across repeated
532 presentations. Future work could try and identify multiple integration periods within the same
533 response by manipulating the type of context which surrounds a segment. Here, we examined
534 two distinct types of contexts and found similar results (**Fig S4**), suggesting that hierarchical
535 temporal integration is a robust property of human auditory cortex.

536
537 Our study focused on characterizing responses within auditory cortex using a relatively short
538 range of segment durations (31 milliseconds to 2 seconds) and a diverse set of natural sounds,
539 including many non-speech and non-music sounds. Future work could characterize integration
540 periods outside of auditory cortex by using a longer range of segment durations and focusing
541 on just speech or music, which exhibit longer-term temporal structure^{8,18,42}. For example, a
542 recent fMRI study provided evidence for multi-second integration periods in regions outside of
543 auditory cortex by examining the delay needed for responses to speech to become context
544 invariant⁵². This study was not able to estimate temporal integration periods within auditory
545 cortex because the timescale of the fMRI response is an order of magnitude slower than
546 auditory cortical integration periods. Here, we took advantage of the rare opportunity to record
547 intracranial responses from the human brain, which is the only human neuroscience method
548 with the spatiotemporal resolution to estimate integration periods in auditory cortex.

549

550 **Methods**

551

552 **Participants & data collection.** Data were collected from 23 patients undergoing treatment
553 for intractable epilepsy at the NYU Langone Hospital (14 patients) and the Columbia University
554 Medical Center (9 patients). One patient was excluded because they had a large portion of the
555 left temporal lobe resected in a prior surgery. Of the remaining 22 subjects, 18 had sound-
556 responsive electrodes (see *Electrode selection*). Electrodes were implanted to localize
557 epileptogenic zones and delineate these zones from eloquent cortical areas before brain
558 resection. NYU patients were implanted with subdural grids, strips and depth electrodes
559 depending on the clinical needs of the patient. CUMC patients were implanted with stereotactic
560 depth electrodes. All subjects gave informed written consent to participate in the study, which
561 was approved by the Institutional Review Boards of CUMC and NYU.

562

563 **Stimuli for TCI paradigm.** Segments were excerpted from 10 natural sound recordings, each
564 two seconds in duration (**Table S1**). Shorter segments were created by subdividing the longer
565 segments. Each natural sound was RMS-normalized before segmentation.

566

567 We tested seven segment durations (31.25, 62.5, 125, 250, 500, 1000, and 2000 ms). We
568 presented the segments of a given duration in two pseudorandom orders, yielding 14
569 sequences (7 durations x 2 orders), each 20 seconds in duration. The only constraint was that
570 a given segment had to be preceded by a different segment in the two orders. When we
571 designed the stimuli, we thought that integration periods might be influenced by transients at
572 the start of a sequence, so we designed the sequences such that the first 2 seconds and last
573 18 seconds of each sequence contained distinct and non-overlapping segments so that we
574 could separately analyze the just last 18 seconds. In practice, integration periods were similar
575 when analyzing the first 18 seconds vs. the entire 20-second sequence. Segments were
576 concatenated using cross-fading to avoid click artifacts (31.25 ms raised cosine window; cross-
577 fading was accounted for in our integration period model). Each stimulus was repeated several
578 times (4 repetitions for most subjects; 8 repetitions for 2 subjects; 6 and 3 repetitions for two
579 other subjects). Stimuli will be made available upon publication.

580

581 **Natural sounds.** In a subset of 11 patients, we measured responses to a diverse set of 119
582 natural sounds from 11 categories, similar to those from our prior studies characterizing
583 auditory cortex¹¹ (there were at least 7 exemplars per category). The sound categories are
584 listed in **Figure 5a**. Most sounds (108) were 4 seconds in duration. The remaining 11 sounds
585 were longer excerpts of English speech (28-70 seconds) that were included to characterize
586 responses to speech for a separate study. Here, we just used responses to the first 4 seconds
587 of these stimuli to make them comparable to the others. The longer excerpts were presented
588 either at the beginning (6 patients) or end of the experiment (5 patients). The non-English
589 speech stimuli were drawn from 10 languages: German, French, Italian, Spanish, Russian,
590 Hindi, Chinese, Swahili, Arabic, Japanese. We classified these stimuli as “foreign speech”
591 since nearly all were unfamiliar to our subjects, though occasionally a patient had some
592 familiarity with one language. Twelve sounds were repeated four times to make it possible to
593 measure response reliability and noise-correct our measures. All other stimuli were presented
594 once. All sounds were RMS normalized.

595

596 As with the main experiment, subjects did not have a formal task but the experiment was
597 periodically paused and subjects were asked a simple question to encourage them to listen to
598 the sounds. For the 4-second sounds, subjects were asked to identify/describe the last sound
599 they heard. For the longer English speech excerpts, subjects were asked to repeat the last
600 phrase they heard.

601

602 **Preprocessing.** Electrode responses were common-average referenced to the grand mean
603 across electrodes from each subject. We excluded noisy electrodes from the common-average
604 reference by detecting anomalies in the 60 Hz power band (measured using an IIR resonance
605 filter with a 3dB down bandwidth of 0.6 Hz). Specifically, we excluded electrodes whose 60 Hz
606 power exceeded 5 standard deviations of the median across electrodes. Because the standard
607 deviation is itself sensitive to outliers, we estimated the standard deviation using the central
608 20% of samples, which are unlikely to be influenced by outliers. Specifically, we divided the
609 range of the central 20% of samples by that which would be expected from a Gaussian of unit
610 variance. After common-average referencing, we used a notch filter to remove harmonics &
611 fractional multiples of the 60 Hz noise (60, 90, 120, 180; using an IIR notch filter with a 3dB
612 down bandwidth of 1 Hz; the filter was applied forward and backward).

613
614 We computed broadband gamma power by measuring the envelope of the preprocessed
615 signal filtered between 70 and 140 Hz (implemented using a 6th order Butterworth filter with
616 3dB down cutoffs of 70 and 140 Hz; the filter was applied forward and backward). Results were
617 similar using other frequency ranges (**Fig S6**). The envelope was measured as the absolute
618 value of the analytic signal after bandpassing. We then downsampled the envelopes to 100
619 Hz (the original sampling rate was either 512, 1000, 1024, or 2048 Hz, depending on the
620 subject). We used simulations to ensure that this procedure would allow us to accurately
621 recover the integration period of broadband power fluctuations (see *Simulations* below).

622
623 Occasionally, we observed visually obvious artifacts in the broadband gamma power for a
624 small number of timepoints. To detect such artifacts, we computed the 90th percentile of each
625 electrode's response distribution across all timepoints. We classified a timepoint as an outlier
626 if it exceeded 5 times the 90th percentile value for each electrode. We found this value to be
627 relatively conservative in that only a small number of timepoints were excluded (on average,
628 0.04% of timepoints were excluded across all sound-responsive voxels). Because there were
629 only a small number of outlier timepoints, we replaced the outlier values with interpolated
630 values from nearby non-outlier timepoints.

631
632 As is standard, we time-locked the iEEG recordings to the stimuli by either cross-correlating
633 the audio with a recording of the audio collected synchronously with the iEEG data or by
634 detecting a series of pulses at the start of each stimulus that were recorded synchronously
635 with the iEEG data. We used the stereo jack on the experimental laptop to either send two
636 copies of the audio or to send audio and pulses on separate channels. The audio on one
637 channel was used to play sounds to subjects, and the audio/pulses on the other were sent to
638 the recording rig. Sounds were played through either a Bose Soundlink Mini II speaker (at
639 CUMC) or an Anker Soundcore speaker (at NYU). Responses were converted to units of
640 percent signal change relative to silence by subtracting and then dividing the response of each
641 electrode by the average response during the 500 ms before each stimulus.

642
643 **Electrode selection.** We selected electrodes with a reliable broadband gamma response to
644 the sound set. Specifically, we measured the test-retest correlation of each electrodes
645 response across all stimuli (using odd vs. even repetitions). We selected electrodes with a
646 test-retest Pearson correlation of at least 0.1, which we found to be sufficient to reliably
647 estimate integration periods in simulations (described below). We ensured that this correlation
648 value was significant using a permutation test, where we randomized the mapping between
649 stimuli across repeated presentations and recomputed the correlation (using 1000
650 permutations). We used a Gaussian fit to the distribution of permuted correlation coefficients
651 to compute small p-values⁵³. Only electrodes with a highly significant correlation relative to the
652 null were kept ($p < 10^{-5}$). We used a low p-value threshold, because we have found that any
653 electrode with a borderline-significant test-retest response across an entire experiment is very

654 noisy. We identified 190 electrodes out of 2847 total that showed a reliable response to natural
655 sounds based on these criteria.

656
657 **Denoising.** Responses in non-primary regions were less reliable on average than responses
658 in primary regions (the median test-retest correlation for each annular ROIs in **Fig 4b** was
659 0.32, 0.23, 0.17). We ensured that differences in reliability could not explain our results in two
660 ways: (1) we ensured that our model-estimated integration periods were unbiased by low data
661 reliability (see *Model-estimated integration periods* and *Simulations* below) (2) we repeated
662 our analyses using a denoising procedure that substantially increased the reliability of the
663 electrode responses to a level well-above the point at which reliability might affect our
664 integration period estimates. Our denoising procedure was motivated by the observation that
665 iEEG responses are relatively low-dimensional and as a consequence much of the stimulus-
666 driven response variation is shared across subjects⁵⁴, in contrast with the noise which differs
667 from subject to subject. We thus projected the electrode responses from one subject onto the
668 responses from all other subjects (using regression), which has the effect of throwing out
669 response variation that is not present in at least two subjects. We have found this procedure
670 to be useful when there are a relatively large number of subjects with responses from a
671 restricted region of the brain like auditory cortex, as was the case in our study. To examine the
672 effect of denoising, we measured split-half reliability before and after denoising (**Fig S10a**).
673 When we denoised both splits of data, the median correlation increased four-fold from 0.21 to
674 0.8 (**Fig S10a**, purple dots). We also found that the reliability improved when only one split
675 was denoised, which indicates that the analysis discarded more noise than signal (reliability
676 improved for 93% of electrodes) (**Fig S10a**, blue dots). Since our results were similar using
677 original and denoised data (compare **Figs 2b&4** with **Figs S10b&c**), we conclude that our
678 findings cannot be explained by differences in data reliability.

679
680 **Electrode localization.** Following standard practice, we localized electrodes as bright spots
681 on a post-operative computer tomography (CT) image or dark spots on a magnetic resonance
682 image (MRI), depending on whichever was available in a given patient. The post-op CT or MRI
683 was aligned to a high-resolution, pre-operative magnetic resonance image (MRI) that was
684 undistorted by electrodes. Each electrode was then projected onto the cortical surface
685 computed by Freesurfer from the pre-op MRI scan, excluding electrodes that were greater than
686 10 mm from the surface. This projection is error prone because faraway points on the cortical
687 surface can be nearby in space due to cortical folding. To minimize gross errors, we
688 preferentially localized sound-responsive electrodes to regions where sound-driven responses
689 are likely to occur⁵⁴. Specifically, we calculated the likelihood of observing a significant
690 response to sound using a recently collected fMRI dataset, where responses were measured
691 to a large set of natural sounds across 20 subjects with whole-brain coverage³³ ($p < 10^{-5}$,
692 measured using a permutation test). We treated this map as a prior and multiplied it by a
693 likelihood map, computed separately for each electrode based on the distance of that electrode
694 to each point on the cortical surface (using a 10 mm FWHM Gaussian error distribution). We
695 then assigned each electrode to the point on the cortical surface where the product of the prior
696 and likelihood was greatest (which can be thought of as the maximum posterior probability
697 solution). We smoothed the prior probability map (10 mm FWHM kernel) so that it would only
698 affect the localization of electrodes at a coarse level, and not bias the location of electrodes
699 locally, and we set the minimum prior probability to be 0.05 to ensure every point had non-zero
700 prior probability. We plot the prior map and its effect on localization in **Fig S11**.

701
702 **Cross-context correlation.** We now review our key analysis for estimating context invariance.
703 MATLAB code implementing these analyses will be made available upon publication. For each
704 electrode and segment duration, we compiled the responses surrounding all segments into a
705 matrix, aligned to segment onset (the segment-aligned response matrix or SIR matrix) (**Fig**

706 **2a).** We calculated a separate SIR matrix for each context, such that corresponding rows
707 contained the response timecourse to the same segment from different contexts. To detect if
708 there was a lag where the response was the same across contexts, we correlated
709 corresponding columns across SIR matrices from different contexts (the cross-context
710 correlation). We compared the cross-context correlation with the correlation when the context
711 was identical using repeated presentations of the same sequence (the noise ceiling).

712
713 The noise ceiling exhibited reliable variation across lags, which is evident from the fact that the
714 cross-context correlation remained yoked to the noise ceiling when the segment duration was
715 long relative to the integration period (evident for example in the HG electrode's data for 250
716 and 500 ms in **Fig 2b**). This variation is expected since the sounds that happen to fall within
717 the integration period will vary with lag, and the noise ceiling will depend upon how strongly
718 the electrode responds to the sounds within the integration period. It is also evident that the
719 variation in the noise ceiling is periodic. This periodicity is an inevitable consequence of
720 correlating across a fixed set of segments. To see this, consider the fact that the onset of one
721 segment is also the offset of the preceding segment. Since we are correlating across segments
722 for a given lag, the values being used to compute the correlation are nearly identical at the
723 start and end of a segment (the only difference occurs for the first and last segment of the
724 entire sequence). The same logic applies to all lags that are separated by a period equal to
725 the segment duration.

726
727 Because the shorter segments were subsets of the longer segments, we could consider two
728 types of context: (1) random context, where a segment is flanked by random other segments
729 (2) natural context, where a segment is a part of a longer segment and thus surrounded by its
730 natural context. Since the two contexts being compared must differ, one of the contexts always
731 has to be random, but the other context can be random or natural. In practice, we found similar
732 results when comparing random-only contexts and when comparing random and natural
733 contexts (**Fig S4**). This fact is practically useful since it greatly increases the number of
734 comparisons that can be made. For example, each 31 millisecond segment had 2 random
735 contexts (one per sequence) and 12 natural contexts (2 sequences x 6 longer segment
736 durations). The two random contexts can be compared with each other as well as with the
737 other 12 natural contexts. For our main analyses, we averaged the cross-context correlation
738 across all of these comparisons for maximal statistical power.

739
740 We note that the cross-context correlation will typically be more reliable for shorter segment
741 durations since there are more segments with which to compute the correlation. We consider
742 this property useful since for electrodes with shorter integration periods there will be a smaller
743 number of lags at the shorter segment durations that effectively determine the integration
744 period, and thus it is useful if these lags are more reliable. Conversely, electrodes with longer
745 integration periods exhibit a more gradual build-up of the cross-context correlation at the longer
746 segment durations, and our model enables us to pool across all of these lags to arrive at a
747 robust estimate of the integration period.

748
749 **Model-estimated integration periods.** We modeled temporal integration periods using a
750 Gamma-distributed window (h) that we scaled and shifted in time:

751

752 (1)
$$h(t; \delta, \lambda, \beta) = g\left(\frac{t - \delta}{\lambda}, \beta\right)$$

753 (2)
$$g(t; \beta) = \frac{\beta^\beta}{\Gamma(\beta)} t^{\beta-1} e^{-\beta t}$$

754

755 The shape is determined by β and varies from more exponential to more bell-shaped (**Fig S5**).
756 The integration width and center do not correspond directly to any of the three parameters
757 (δ, λ, β), mainly because the scale parameter (λ) alters both the center and width. The
758 integration width was defined as the smallest interval that contained 75% of the window's
759 mass, and the integration center was defined as the window's median. Both parameters were
760 calculated numerically from the cumulative distribution function of the window.

761
762 For a given integration window, we predicted the cross-context correlation at each lag and
763 segment duration by measuring how much the integration period overlaps the shared central
764 segment (w) vs. the N surrounding context segments (β_n) (see **Fig 3b**):
765

766 (3)
$$r_{ceil} \frac{w^2}{w^2 + \sum_{n=1}^N \beta_n^2}$$

767
768 where r_{ceil} is the measured noise ceiling, and the ratio on the right is the predicted correlation
769 in the absence of noise. The predicted cross-context correlation varies with the segment
770 duration and lag because the overlap varies with the segment duration and lag. When the
771 integration period only overlaps the shared segment ($w = 1, \sum \beta_n = 0$), the model predicts a
772 correlation of 1 in the absence of noise, and when the integration period only overlaps the
773 surrounding context segments ($w = 0, \sum \beta_n = 1$), the model predicts a correlation of 0. In
774 between these two extremes, the predicted cross-context correlation equals the fraction of the
775 response driven by the shared segment, with the response variance for each segment given
776 by the squared overlap with the integration period. A formal derivation of this equation is given
777 at the end of the Methods (see *Deriving a prediction for the cross-context correlation*). For a
778 given segment duration, the overlap with each segment was computed by convolving the
779 model integration period with a boxcar function whose width is equal to the segment duration
780 (with edges tapered to account for cross-fading).

781
782 We varied the width, center and shape of the model integration period and selected the window
783 with the smallest prediction error. Since the cross-context correlation is more reliable for
784 shorter segment durations due to the greater number of segments, we weighted the error by
785 the number of segments used to compute the correlation before averaging across segment
786 durations. Integration widths varied between 31.25 and 1 second (using 100 logarithmically
787 spaced steps). Integration centers varied from the minimum possible given for a causal window
788 up to 500 milliseconds beyond the minimum in 10 millisecond steps. We tested five window
789 shapes ($\beta = 1, 2, 3, 4, 5$).

790
791 We found in simulations that there was an upward bias in the estimated integration widths for
792 noisy data when using the mean squared error (see *Simulations* below). We checked that this
793 bias did not affect our results in two ways. First, we repeated our analyses using denoised
794 data, whose reliability was well above the point at which the bias has any effect (**Fig S10**).
795 Second, we derived a bias-corrected metric, which substantially reduced the bias in
796 simulations (see *Bias correction* below). We used this bias-corrected metric for all of our
797 analyses, but found that the results were very similar with and without bias correction,
798 indicating that our data were sufficiently reliable to avoid any substantial bias even without
799 denoising (compare **Fig 4** which shows results with correction with **Fig S12** which shows
800 results without correction).

801
802 We assessed the significance of our model predictions by creating a null distribution using
803 phase-scrambled model predictions. Phase scrambling exactly preserves the mean, variance
804 and autocorrelation of the predictions but alters the locations of the peaks and valleys. Phase
805 scrambling was implemented by shuffling the phases of different frequency components

806 without altering their amplitude and then reconstructing the signal (using the FFT/iFFT). After
807 phase-scrambling, we remeasured the error between the predicted and measured cross-
808 context correlation, and selected the model with the smallest error (as was done for the
809 unscrambled predictions). We repeated this procedure 100 times to build up a null distribution,
810 and used this null distribution to calculate a p-value for the actual error based on unscrambled
811 predictions (again fitting the null distribution with a Gaussian to calculate small p-values). For
812 95% of sound-responsive electrodes (181 of 190), the model's predictions were highly
813 significant ($p < 10^{-5}$).

814
815 **Simulations.** We tested our complete analysis pipeline using simulated data. Specifically, we
816 modulated a broadband carrier (Gaussian noise filtered between 70 and 140 Hz) with the
817 waveform amplitude of the stimuli from our TCI paradigm, integrated within a Gamma-
818 distributed integration period. We added Gaussian noise to manipulate the test-retest reliability
819 of these responses and thus determine the minimum reliability needed to accurately infer
820 integration periods. We simulated four different responses to the same stimuli using
821 independent samples of the broadband carrier and additive noise (for most subjects we had
822 four repetitions), and iteratively increased or decreased the noise level to achieve a desired
823 split-half correlation (within a tolerance of 0.001).

824
825 We then applied our complete analysis pipeline to these simulated responses. Our goal was
826 to assess if we could accurately infer the true integration width from the simulated data. We
827 thus varied the integration width of the simulated response (from 31 ms to 500 ms in octave
828 steps), using a fixed shape ($\beta = 3$) and center (set to the minimum value for a casual window).
829 However, we did not assume that the shape or center were known, and thus varied the shape
830 and center along with the width when inferring the best-fit integration period, as was done for
831 the iEEG analyses.

832
833 We found that the estimated integration widths were close to the true widths when the split-
834 half reliability was at least 0.1 (**Fig S13**), which was the reliability cutoff used to select sound-
835 responsive electrodes. When the test-retest reliability was low, there was an upwards bias in
836 the estimated widths when using the mean squarer error (**Fig S13**, top panel), which was
837 substantially reduced using our bias-corrected metric (**Fig S13**, bottom panel; see *Bias*
838 *correction* below).

839
840 **Anatomical ROI analyses.** We grouped electrodes into regions-of-interest (ROI) based on
841 their anatomical distance to posteromedial Heschl's gyrus (TE1.1)⁵⁵ (**Fig 4b**), which is a
842 common anatomical landmark for primary auditory cortex^{31,56}. Distance was measured on the
843 flattened 2D representation of the cortical surface as computed by Freesurfer. Electrodes were
844 grouped into three 10 millimeter bins (0-10, 10-20, and 20-30 mm), and we measured the
845 median integration width and center across the electrodes in each bin, separately for each of
846 the two hemispheres.

847
848 Error bars and significance were computed by bootstrapping across subjects. Specifically, we
849 resampled subjects with replacement 10,000 times and recomputed the median integration
850 width or center within each ROI using the resampled dataset. A small fraction of these samples
851 (1.2%) were discarded because the resampled dataset did not contain any electrodes in one
852 of the six ROIs (3 distances x 2 hemispheres). To assess whether two ROIs significantly
853 differed (e.g. nearest vs. farthest, left vs. right), we counted the fraction of samples where the
854 resampled values were consistently higher or lower in one of the two ROIs (whichever fraction
855 was lower), subtracted this fraction from 1, and multiplied by 2 to arrive at a two-sided p-value.
856

857 **Category-selective components at different temporal integration periods.** To investigate
858 selectivity for categories, we used responses to the larger set of 119 natural sounds that were
859 tested in a subset of 11 patients. There were 104 electrodes from these 11 subjects that
860 passed the inclusion criteria described above (out of 181 total). We grouped these electrodes
861 based on the width of their integration period in octave intervals, spaced a half-octave apart
862 (intervals shown in **Fig 5a**). Each group had between 22 and 43 electrodes. We then used
863 several different analyses to investigate the degree of category selectivity in each group.

864
865 We used a combination of component methods (**Fig 5a,c**) and individual-electrode analyses
866 (**Fig 5b**) to assess category selectivity. Component methods are commonly used to summarize
867 responses from a population of electrodes or neurons⁵⁷. And we have previously shown that
868 component methods can better isolate selectivity for categories compared with analyzing
869 individual iEEG electrodes⁵⁴ or individual fMRI voxels¹¹. To visualize the dominant structure at
870 each integration period, we projected the responses onto the top two principle components
871 (PCs) from each group of electrodes (**Fig S8**). If these components exhibit selectivity for
872 categories then the average component response to different categories should appear
873 segregated when plotted as a trajectory. Because the first two PCs might obscure category
874 selectivity present at higher PCs, we repeated the analysis using the two components that best
875 separated the categories, estimated using linear discriminant analysis (LDA)⁵⁸ (**Fig 5a**). LDA
876 was applied to timepoints between 250 milliseconds and 4 seconds after stimulus onset to
877 account for response delays. To avoid statistical circularity, we used half the sounds to infer
878 components, and the other half to measure their response. And to prevent the analysis from
879 targeting extremely low-variance components, we applied LDA to the top five PCs from each
880 electrode group.

881
882 PCs were computed using responses from the TCI experiment, where we had responses from
883 a larger number of electrodes and subjects (181 electrodes from 18 subjects). We then
884 estimated the response of these same PCs to the larger set of 119 natural sounds using the
885 subset of electrodes with responses in both experiments (104 electrodes from 11 subjects).
886 Since each PC is just a weighted sum of the electrode responses, we simply multiplied the
887 responses to the 119 natural sounds by the reconstruction weights inferred from the TCI
888 experiment. Since only a subset of electrodes were tested in both experiments, we inferred
889 the reconstruction weights using just the electrodes tested in both experiments, by finding the
890 linear combination of these electrodes that best approximated each PC.

891
892 **Feature predictions.** As a complement to the component analyses, we measured the degree
893 to which individual electrode responses could be predicted from category labels (**Fig 5b**). We
894 binned the results based on integration width of the electrode, using the same octave-spaced
895 intervals. And we compared the prediction accuracies for the category labels with those from
896 a cochleagram representation of sound.

897
898 Cochleagrams were calculated using a cosine filterbank with bandwidths designed to mimic
899 cochlear tuning³¹ (29 filters between 50 Hz and 20 kHz, 2x overcomplete). The envelopes from
900 the output of each filter were compressed to mimic cochlear amplification (0.3 power). The
901 frequency axis was resampled to a resolution of 12 cycles per octave and the time axis was
902 resampled to 100 Hz (the sampling rate used for all of our analyses).

903
904 For each category label, we created a binary timecourse with 1s for all timepoints/sounds from
905 that category, and 0s for all other timepoints. We only labeled timepoints with a 1 if they had
906 sound energy that exceeded a minimum threshold. The sound energy at each moment in time
907 was calculated by averaging the cochleagram across frequency, and the minimum threshold
908 was set to one fifth the mean energy across all timepoints and sounds.

909

910 We predicted electrode responses between 500 milliseconds pre-stimulus onset to 4 seconds
911 post-stimulus onset. We used ridge regression to learn a linear mapping between these
912 features and the response. We included five delayed copies of each regressor, with the delays
913 selected to span the integration period of the electrode (from the bottom fifth to the top fifth
914 quintile). Regression weights were fit using the 107 sounds that were presented once, and we
915 evaluated the fits using the 12 test sounds that were repeated four times each, making it
916 possible to compute a noise-corrected measure of prediction accuracy^{59,60}:

917

$$918 \quad (4) \quad \frac{[0.5 * \text{corr}(r_1, p) + 0.5 * \text{corr}(r_2, p)]^2}{\text{corr}(r_1, r_2)}$$

919

920 where r_1 and r_2 are two independent measures of the response (computed using odd and even
921 repetitions) and p is the prediction computed from the training data. We used cross-validation
922 within the training set to choose the regularization coefficient (testing a wide range of values
923 from 2^{-100} to 2^{100} in octave steps).

924

925 Significance and error bars were calculated using bootstrapping across subjects (the same
926 procedure described above in *Anatomical ROI analyses*). To test whether the prediction
927 accuracy increased or decreased as a function of the integration width, we measured the slope
928 between the noise-corrected prediction accuracy and the integration width (on a logarithmic
929 scale). We then tested whether the bootstrapped slopes for the category and cochlear
930 predictions differed significantly from zero and from each other.

931

932 **Speech and music selectivity.** We separately quantified the degree of speech and music
933 selectivity at each integration width (**Fig 5c**). Selectivity was quantified as the degree of
934 separation between speech/music and all other sounds, along the components that showed
935 the greatest selectivity for speech/music. Both English and foreign speech sounds were
936 grouped as speech, since they yielded similar responses, consistent with prior results^{11,12}.
937 Vocal music was excluded from the speech selectivity analysis since it produced an
938 intermediate responses along the speech-selective component (**Fig S9**), as expected since
939 vocals contain speech¹¹. Instrumental music, vocal music and drumming were grouped as
940 music, since they produce above average responses in music-selective brain regions³³.

941

942 The weights for each component were learned by regressing the electrode responses from
943 each group against a binary category vector with 1s for all timepoints/sounds from the target
944 category (e.g. speech) and 0s for all other timepoints/sounds. To avoid over-fitting to low-
945 variance signals, we again applied our analysis to the top 5 PCs from each group. We used
946 independent sounds to estimate components and measure their response (5-fold cross-
947 validation, each fold had a similar number of sounds per category). To account for the
948 response delay we only used responses between 250 milliseconds and 4 seconds post-
949 stimulus onset.

950

951 **Figure S9** plots the average component response timecourse to each sound category
952 (averaged across test folds). **Figure 5c** plots the average separation, measured as the d-prime
953 between responses to sounds from the target and non-target category across all timepoints
954 (between 250 milliseconds and 4 seconds post-stimulus onset). D-prime is a standard
955 measure of the separation between two responses, and is defined as the difference in the
956 mean response divided by square root of the average variances:

957

958 (5)

$$\frac{\mu_1 - \mu_2}{\sqrt{\frac{\sigma_1^2}{2} + \frac{\sigma_2^2}{2}}}$$

959

960 The average within-category variance in the denominator of equation 5 is a sum of the
961 stimulus-driven response variance, which is repeatable across measurements, and the noise
962 variance, which is not (the means are unbiased by noise assuming the noise is zero mean).
963 We therefore noise-corrected our d-prime measure by subtracting off an estimate of the noise
964 variance from the measured within-category variance. We estimated the noise variance as half
965 of the error variance to repeated presentations of the same stimulus (using the 12 sounds
966 repeated 4 times)³¹. We used half of the error variance, since the error reflects the difference
967 between two independent measurements of the same signal, and the total variance of two
968 independent signals that are subtracted or summed is additive.

969

970 Significance and error bars were computed via bootstrapping across sounds. Unlike other
971 analyses, we did not bootstrap across subjects because doing so would have been
972 inappropriate for this particular analysis. Specifically, each component was computed by
973 regressing the electrode responses from all subjects against the target category vector, and
974 thus in the language of regression, electrodes/subjects are features and timepoints/sounds
975 are observations. While bootstrapping across observations (timepoints/sounds) is standard⁶¹,
976 bootstrapping across features (electrodes/subjects) is inappropriate, because repeating
977 features does not change the least-squares solution.

978

979 We bootstrapped across sounds by sampling sounds with replacement, separately for the
980 target and non-target category (e.g. speech and non-speech sounds). We also resampled
981 sounds from the test set used to calculate the noise variance. We then recalculated the noise-
982 corrected d-prime using the component timecourses for the resampled sounds (repeating
983 timecourses for sounds sampled more than once). To test if selectivity increased with the
984 integration width, we measured if the bootstrapped slope between selectivity and the
985 integration width was significantly greater than 0.

986

987 **Deriving a prediction for the cross-context correlation.** We now derive the equation used
988 to predict the cross-context correlation from a model integration period (equation 3). The cross-
989 context correlation is computed across segments for a fixed lag and segment duration by
990 correlating corresponding columns of SIR matrices from different contexts (**Fig 2a**). Consider
991 two pairs of cells ($e_{s,A}$, $e_{s,B}$) from these SIR matrices, representing the response to a single
992 segment (s) in two different contexts (A, B) for a fixed lag and segment duration (we do not
993 indicate the lag and segment duration to simplify notation). To reason about how the shared
994 and context segments might relate to the cross-context correlation at each moment in time,
995 we assume that the response reflects the sum of the responses to each segment weighted by
996 the degree of overlap with the integration period (**Fig 3b**):

997

998 (6)

$$e_{s,A} = wr(s) + \sum_{n=1}^N \beta_n r(c_{s,A,n})$$

999 (7)

$$e_{s,B} = wr(s) + \sum_{n=1}^N \beta_n r(c_{s,B,n})$$

1000

1001 where $r(s)$ reflects the response to the shared central segment, $r(c_{s,A,n})$ and $r(c_{s,B,n})$ reflect
1002 the response to the n-th surrounding segment in each of the two contexts (i.e. the segment

1003 right before and right after, two before and two after, etc.), and w and β_n reflect the degree of
 1004 overlap with the shared and surrounding segments, respectively (illustrated in **Fig 3b**).

1006 Below we write down the expectation of the cross-context correlation in the absence of noise,
 1007 substitute equations 6 & 7, and simplify. Moving from line 8 to line 9 takes advantage of the
 1008 fact that contexts A and B are no different in structure and so their expected variance is the
 1009 same. Moving from line 10 to line 11, we have taken advantage of the fact that surrounding
 1010 context segments are random, and thus all cross products that involve the context segments
 1011 are zero in expectation, canceling out all of the terms except those noted in equation 11.
 1012 Finally, in moving from equation 11 to 12, we take advantage of the fact that there is nothing
 1013 special about the segments that make up the shared central segments compared with the
 1014 surrounding context segments, and their expected variance is therefore equal and cancels
 1015 between the numerator and denominator.

$$1017 \quad (8) \quad E[r_{cross}] = \frac{E_s[e_{s,A}e_{s,B}]}{\sqrt{E_s[e_{s,A}^2]E_s[e_{s,B}^2]}}$$

$$1018 \quad (9) \quad = \frac{E_s[e_{s,A}e_{s,B}]}{E_s[e_{s,A}^2]}$$

$$1019 \quad (10) \quad = \frac{E_s[(wr(s) + \sum_{n=1}^N \beta_n r(c_{s,A,n}))(wr(s) + \sum_{n=1}^N \beta_n r(c_{s,B,n}))]}{E_s[(wr(s) + \sum_{n=1}^N \beta_n r(c_{s,A,n}))^2]}$$

$$1020 \quad (11) \quad = \frac{w^2 E_s[r^2(s)]}{w^2 E_s[r^2(s)] + \sum_{n=1}^N \beta_n^2 E_s[r(c_{s,A,n})^2]}$$

$$1021 \quad (12) \quad = \frac{w^2}{w^2 + \sum_{n=1}^N \beta_n^2}$$

1022
 1023 We multiplied equation 12 by the noise ceiling to arrive at our prediction of the cross context
 1024 correlation (equation 3).

1026 **Bias correction.** Here, we derive the correction procedure used to minimize the bias when
 1027 evaluating model predictions via the squared error.

1029 Before beginning, we highlight a potentially confusing, but necessary distinction between noisy
 1030 measures and noisy data. As we show below, the bias is caused by the fact that our correlation
 1031 measures are noisy in the sense that they will not be the same across repetitions of the
 1032 experiment. The bias is *not* directly caused by the fact that the data is noisy, since if there are
 1033 enough segments the correlation measures will be reliable even if the data are noisy, which is
 1034 what matters since we explicitly measure and account for the noise ceiling. To avoid confusion,
 1035 we use the superscript (n) to indicate noisy measures, (t) to indicate the true value of a noisy
 1036 measure (i.e. in the limit of infinite segments), and (p) to indicate a “pure” measure computed
 1037 from noise-free data.

1039 Consider the error between the measured $(r_{cross}^{(n)})$ and model-predicted $(p_{cross}^{(n)})$ cross-context
 1040 correlation for a single lag and segment duration (the model prediction is noisy because of
 1041 multiplication with the noise ceiling which is measured from data):

1042

1043 (13)
$$\left[r_{cross}^{(n)} - p_{cross}^{(n)} \right]^2$$

1044
 1045 Our final cost function averaged these pointwise squared errors across all lags and segment
 1046 durations weighted by the number of segments used to compute each correlation (which was
 1047 greater for shorter segment durations). Here, we analyze each lag and segment duration
 1048 separately, and thus ignore the influence of the weights which is simply a multiplicative factor
 1049 that can be applied at the end after bias correction.

1050
 1051 Our analysis proceeds by writing the measured ($r_{cross}^{(n)}$) and predicted ($p_{cross}^{(n)}$) cross-context
 1052 correlation in terms of their underlying true and pure measures (equations 14 to 17). We then
 1053 substituting these definitions into the expectation of the squared error and simplify (equations
 1054 18 to 21), which yields insight into the cause of the bias.

1056 The cross-context correlation ($r_{cross}^{(n)}$) is the sum of the true cross-context correlation plus error:

1057
 1058 (14)
$$r_{cross}^{(n)} = r_{cross}^{(t)} + e_{cross}$$

1059
 1060 And the true cross-context correlation is the product of the pure/noise-free cross-context
 1061 correlation ($r_{cross}^{(p)}$) with the true noise ceiling ($r_{ceil}^{(t)}$):

1062
 1063 (15)
$$r_{cross}^{(t)} = r_{cross}^{(p)} r_{ceil}^{(t)}$$

1064
 1065 The predicted cross-context correlation is the product of the noise-free prediction ($p_{cross}^{(p)}$) times
 1066 the measured noise ceiling ($r_{ceil}^{(n)}$):

1067
 1068 (16)
$$p_{cross}^{(n)} = p_{cross}^{(p)} r_{ceil}^{(n)}$$

1069
 1070 And the measured noise ceiling is the sum of the true noise ceiling ($r_{ceil}^{(t)}$) plus error ($e_{ceil}^{(n)}$):

1071
 1072 (17)
$$r_{ceil}^{(n)} = r_{ceil}^{(t)} + e_{ceil}$$

1073
 1074 Below we substitute the above equations into the expectation for the squared error and
 1075 simplify. Only the error terms (e_{cross} and e_{ceil}) are random variables, and thus in equation 20,
 1076 we have moved all of the other terms out of the expectation. In moving from equations 20 to
 1077 21, we make the assumption / approximation that the errors are uncorrelated and zero mean,
 1078 which causes all but three terms to dropout in equation 21. This approximation, while possibly
 1079 imperfect, substantially simplifies the expectation and makes it possible to derive a simple bias
 1080 correction procedure, as described next.

1081
 1082 (18)
$$E \left[\left(r_{cross}^{(n)} - p_{cross}^{(n)} \right)^2 \right] = E \left[\left(\left(r_{cross}^{(p)} r_{ceil}^{(t)} + e_{cross} \right) - \left(p_{cross}^{(p)} \left(r_{ceil}^{(t)} + e_{ceil} \right) \right) \right)^2 \right]$$

1083 (19)
$$= E \left[\left(r_{ceil}^{(t)} \left(r_{cross}^{(p)} - p_{cross}^{(p)} \right) + e_{cross} - p_{cross}^{(p)} e_{ceil} \right)^2 \right]$$

1084 (20)
$$= r_{ceil}^{(t)2} \left(r_{cross}^{(p)} - p_{cross}^{(p)} \right)^2 + E[e_{cross}^2] + \left(p_{cross}^{(p)} \right)^2 E[e_{ceil}^2]$$

1085
$$+ 2r_{ceil}^{(t)} \left(r_{cross}^{(p)} - p_{cross}^{(p)} \right) E[e_{cross}] - 2r_{ceil}^{(t)} \left(r_{cross}^{(p)} - p_{cross}^{(p)} \right) p_{cross}^{(p)} E[e_{ceil}]$$

$$\begin{aligned} & -2p_{cross}^{(p)} E[e_{ceil}e_{cross}] \\ (21) \quad & \approx r_{ceil}^{(t)2} \left(r_{cross}^{(p)} - p_{cross}^{(p)} \right)^2 + E[e_{cross}^2] + \left(p_{cross}^{(p)} \right)^2 E[e_{ceil}^2] \end{aligned}$$

The first term in equation 21 is what we would hope to measure: a factor which is proportional to the squared error between the pure cross-context correlation computed from noise-free data ($r_{cross}^{(p)}$) and the model's prediction of the pure cross-context correlation ($p_{cross}^{(p)}$). The second term does not depend upon the model's prediction and thus can be viewed as a constant from the standpoint of analyzing model bias. The third term is potentially problematic, since it biases the error upwards based on the squared magnitude of the predictions, with the magnitude of the bias determined by the magnitude of the errors in the noise ceiling. This term results in an upward bias in the estimated integration width, because narrower integration periods have larger squared magnitudes on average. This bias is only present when there is substantial error in the noise ceiling, which explains why we only observed the bias for data with low reliability (**Fig S13**, top panel).

We can correct for this bias by subtracting a factor whose expectation is equal to the problematic third term in equation 21. All we need is a sample of the error in the noise ceiling, which our procedure naturally provides since we measure the noise ceiling separately for segments from each of the two contexts and then average these two estimates. Thus, we can get a sample of the error by subtracting our two samples of the correlation ceiling and dividing by 2 (averaging is equivalent to summing and dividing by 2 and the noise power of summed and subtracted signals is equal). We then take this sample of the error multiply it by our model prediction, square the result, and subtract this number from the measured squared error. This procedure is done separately for every lag and segment duration.

We found this procedure substantially reduced the bias when pooling across both random and natural contexts (**Fig S13**, bottom panel), as was done for all of our analyses except those shown in **Figure S4**. When only considering random contexts, we found this procedure somewhat over-corrected the bias (inducing a downward bias for noisy data), perhaps due to the influence of the terms omitted in our approximation (equation 21). However, our results were very similar when using random or natural contexts (**Fig S4**), when using either the uncorrected (**Fig S12**) or bias-corrected error (**Fig 2**), and when using highly denoised data (**Fig S10**). Thus, we conclude that our findings were not substantially influenced by noise and were robust to details of the analysis.

References

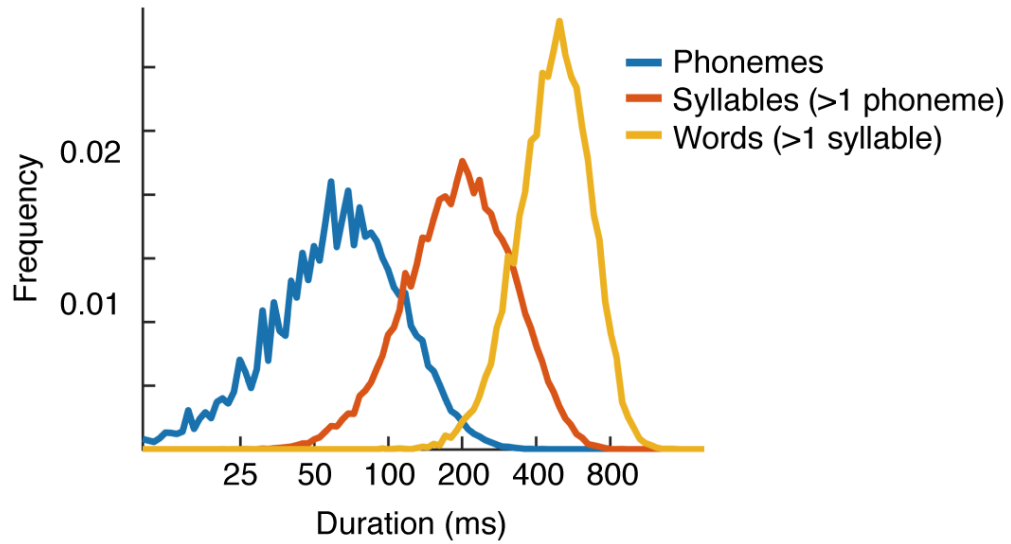
- 1120
1121
1122 1. Chomsky, N. & Halle, M. The sound pattern of English. (1968).
1123 2. Zatorre, R. J., Belin, P. & Penhune, V. B. Structure and function of auditory cortex: music
1124 and speech. *Trends in Cognitive Sciences* **6**, 37–46 (2002).
1125 3. Hickok, G. & Poeppel, D. The cortical organization of speech processing. *Nature reviews*
1126 *neuroscience* **8**, 393–402 (2007).
1127 4. Rauschecker, J. P. & Scott, S. K. Maps and streams in the auditory cortex: nonhuman
1128 primates illuminate human speech processing. *Nature Neuroscience* **12**, 718–724 (2009).
1129 5. Poeppel, D. The analysis of speech in different temporal integration windows: cerebral
1130 lateralization as ‘asymmetric sampling in time’. *Speech communication* **41**, 245–255
1131 (2003).
1132 6. Albouy, P., Benjamin, L., Morillon, B. & Zatorre, R. J. Distinct sensitivity to
1133 spectrotemporal modulation supports brain asymmetry for speech and melody. *Science*
1134 **367**, 1043–1047 (2020).
1135 7. Sharpee, T. O., Atencio, C. A. & Schreiner, C. E. Hierarchical representations in the
1136 auditory cortex. *Current opinion in neurobiology* **21**, 761–767 (2011).
1137 8. Ding, N., Melloni, L., Zhang, H., Tian, X. & Poeppel, D. Cortical tracking of hierarchical
1138 linguistic structures in connected speech. *Nature Neuroscience* **19**, 158–164 (2016).
1139 9. Belin, P., Zatorre, R. J., Lafaille, P., Ahad, P. & Pike, B. Voice-selective areas in human
1140 auditory cortex. *Nature* **403**, 309–312 (2000).
1141 10. Leaver, A. M. & Rauschecker, J. P. Cortical representation of natural complex sounds:
1142 effects of acoustic features and auditory object category. *The Journal of Neuroscience*
1143 **30**, 7604–7612 (2010).
1144 11. Norman-Haignere, S. V., Kanwisher, N. G. & McDermott, J. H. Distinct cortical pathways
1145 for music and speech revealed by hypothesis-free voxel decomposition. *Neuron* **88**,
1146 1281–1296 (2015).
1147 12. Overath, T., McDermott, J. H., Zarate, J. M. & Poeppel, D. The cortical analysis of speech-
1148 specific temporal structure revealed by responses to sound quilts. *Nature neuroscience*
1149 **18**, 903–911 (2015).
1150 13. Davis, M. H. & Johnsruide, I. S. Hierarchical processing in spoken language
1151 comprehension. *Journal of Neuroscience* **23**, 3423–3431 (2003).
1152 14. Di Liberto, G. M., O’Sullivan, J. A. & Lalor, E. C. Low-frequency cortical entrainment to
1153 speech reflects phoneme-level processing. *Current Biology* **25**, 2457–2465 (2015).
1154 15. Mesgarani, N., Cheung, C., Johnson, K. & Chang, E. F. Phonetic feature encoding in
1155 human superior temporal gyrus. *Science* **343**, 1006–1010 (2014).
1156 16. Steinschneider, M., Nourski, K. V. & Fishman, Y. I. Representation of speech in human
1157 auditory cortex: is it special? *Hearing research* **305**, 57–73 (2013).
1158 17. Theunissen, F. & Miller, J. P. Temporal encoding in nervous systems: A rigorous
1159 definition. *J Comput Neurosci* **2**, 149–162 (1995).
1160 18. Lerner, Y., Honey, C. J., Silbert, L. J. & Hasson, U. Topographic mapping of a hierarchy
1161 of temporal receptive windows using a narrated story. *Journal of Neuroscience* **31**, 2906–
1162 2915 (2011).
1163 19. Hullett, P. W., Hamilton, L. S., Mesgarani, N., Schreiner, C. E. & Chang, E. F. Human
1164 superior temporal gyrus organization of spectrotemporal modulation tuning derived from
1165 speech stimuli. *Journal of Neuroscience* **36**, 2014–2026 (2016).
1166 20. Meyer, A. F., Williamson, R. S., Linden, J. F. & Sahani, M. Models of neuronal stimulus-
1167 response functions: elaboration, estimation, and evaluation. *Frontiers in systems*
1168 *neuroscience* **10**, 109 (2017).
1169 21. Barton, B., Venezia, J. H., Saberi, K., Hickok, G. & Brewer, A. A. Orthogonal acoustic
1170 dimensions define auditory field maps in human cortex. *Proceedings of the National*
1171 *Academy of Sciences* **109**, 20738–20743 (2012).

- 1172 22. Harper, N. S. *et al.* Network receptive field modeling reveals extensive integration and
1173 multi-feature selectivity in auditory cortical neurons. *PLoS computational biology* **12**,
1174 e1005113 (2016).
- 1175 23. Flinker, A., Doyle, W. K., Mehta, A. D., Devinsky, O. & Poeppel, D. Spectrotemporal
1176 modulation provides a unifying framework for auditory cortical asymmetries. *Nature*
1177 *human behaviour* **3**, 393 (2019).
- 1178 24. Schönwiesner, M. & Zatorre, R. J. Spectro-temporal modulation transfer function of single
1179 voxels in the human auditory cortex measured with high-resolution fMRI. *Proceedings of*
1180 *the National Academy of Sciences* **106**, 14611–14616 (2009).
- 1181 25. Rogalsky, C., Rong, F., Saberi, K. & Hickok, G. Functional anatomy of language and
1182 music perception: temporal and structural factors investigated using functional magnetic
1183 resonance imaging. *The Journal of Neuroscience* **31**, 3843–3852 (2011).
- 1184 26. Angeloni, C. & Geffen, M. N. Contextual modulation of sound processing in the auditory
1185 cortex. *Current opinion in neurobiology* **49**, 8–15 (2018).
- 1186 27. Griffiths, T. D. *et al.* Direct recordings of pitch responses from human auditory cortex.
1187 *Current Biology* **20**, 1128–1132 (2010).
- 1188 28. Ray, S., Crone, N. E., Niebur, E., Franaszczuk, P. J. & Hsiao, S. S. Neural correlates of
1189 high-gamma oscillations (60–200 Hz) in macaque local field potentials and their potential
1190 implications in electrocorticography. *Journal of Neuroscience* **28**, 11526–11536 (2008).
- 1191 29. Steinschneider, M., Fishman, Y. I. & Arezzo, J. C. Spectrotemporal analysis of evoked
1192 and induced electroencephalographic responses in primary auditory cortex (A1) of the
1193 awake monkey. *Cerebral Cortex* **18**, 610–625 (2008).
- 1194 30. Slaney, M. Auditory toolbox. *Interval Research Corporation, Tech. Rep* **10**, 1998 (1998).
- 1195 31. Norman-Haignere, S. V. & McDermott, J. H. Neural responses to natural and model-
1196 matched stimuli reveal distinct computations in primary and nonprimary auditory cortex.
1197 *PLoS biology* **16**, e2005127 (2018).
- 1198 32. Brodbeck, C., Hong, L. E. & Simon, J. Z. Rapid transformation from auditory to linguistic
1199 representations of continuous speech. *Current Biology* **28**, 3976–3983 (2018).
- 1200 33. Boebinger, D., Norman-Haignere, S., McDermott, J. & Kanwisher, N. Cortical music
1201 selectivity does not require musical training. *bioRxiv* (2020).
- 1202 34. Joris, P. X., Schreiner, C. E. & Rees, A. Neural processing of amplitude-modulated
1203 sounds. *Physiological reviews* **84**, 541–577 (2004).
- 1204 35. Wang, X., Lu, T., Bendor, D. & Bartlett, E. Neural coding of temporal information in
1205 auditory thalamus and cortex. *Neuroscience* **154**, 294–303 (2008).
- 1206 36. Atiani, S. *et al.* Emergent selectivity for task-relevant stimuli in higher-order auditory
1207 cortex. *Neuron* **82**, 486–499 (2014).
- 1208 37. Mizrahi, A., Shalev, A. & Nelken, I. Single neuron and population coding of natural sounds
1209 in auditory cortex. *Current opinion in neurobiology* **24**, 103–110 (2014).
- 1210 38. Leonard, M. K., Bouchard, K. E., Tang, C. & Chang, E. F. Dynamic encoding of speech
1211 sequence probability in human temporal cortex. *Journal of Neuroscience* **35**, 7203–7214
1212 (2015).
- 1213 39. Di Liberto, G. M., Wong, D., Melnik, G. A. & de Cheveigné, A. Low-frequency cortical
1214 responses to natural speech reflect probabilistic phonotactics. *Neuroimage* **196**, 237–247
1215 (2019).
- 1216 40. Gattass, R., Gross, C. G. & Sandell, J. H. Visual topography of V2 in the macaque. *Journal*
1217 *of Comparative Neurology* **201**, 519–539 (1981).
- 1218 41. Dumoulin, S. O. & Wandell, B. A. Population receptive field estimates in human visual
1219 cortex. *Neuroimage* **39**, 647–660 (2008).
- 1220 42. Patel, A. D. *Music, language, and the brain.* (Oxford university press, 2007).
- 1221 43. Elhilali, M. Modulation representations for speech and music. in *Timbre: Acoustics,*
1222 *Perception, and Cognition* 335–359 (Springer, 2019).

- 1223 44. Henry, M. J., Herrmann, B. & Obleser, J. Selective attention to temporal features on
1224 nested time scales. *Cerebral Cortex* **25**, 450–459 (2015).
- 1225 45. Panzeri, S., Brunel, N., Logothetis, N. K. & Kayser, C. Sensory neural codes using
1226 multiplexed temporal scales. *Trends in Neurosciences* **33**, 111–120 (2010).
- 1227 46. Walker, K. M., Bizley, J. K., King, A. J. & Schnupp, J. W. Multiplexed and robust
1228 representations of sound features in auditory cortex. *Journal of Neuroscience* **31**, 14565–
1229 14576 (2011).
- 1230 47. Osman, A. F., Lee, C. M., Escabí, M. A. & Read, H. L. A hierarchy of time scales for
1231 discriminating and classifying the temporal shape of sound in three auditory cortical fields.
1232 *Journal of Neuroscience* **38**, 6967–6982 (2018).
- 1233 48. Ulanovsky, N., Las, L., Farkas, D. & Nelken, I. Multiple time scales of adaptation in
1234 auditory cortex neurons. *Journal of Neuroscience* **24**, 10440–10453 (2004).
- 1235 49. Lu, K. *et al.* Implicit memory for complex sounds in higher auditory cortex of the ferret.
1236 *Journal of Neuroscience* **38**, 9955–9966 (2018).
- 1237 50. Chew, S. J., Mello, C., Nottebohm, F., Jarvis, E. & Vicario, D. S. Decrements in auditory
1238 responses to a repeated conspecific song are long-lasting and require two periods of
1239 protein synthesis in the songbird forebrain. *Proceedings of the National Academy of*
1240 *Sciences* **92**, 3406–3410 (1995).
- 1241 51. Bianco, R. *et al.* Long-term implicit memory for sequential auditory patterns in humans.
1242 *eLife* **9**, e56073 (2020).
- 1243 52. Chien, H.-Y. S. & Honey, C. J. Constructing and Forgetting Temporal Context in the
1244 Human Cerebral Cortex. *Neuron* (2020).
- 1245 53. Norman-Haignere, S. V. *et al.* Pitch-responsive cortical regions in congenital amusia. *J.*
1246 *Neurosci.* **36**, 2986–2994 (2016).
- 1247 54. Norman-Haignere, S. V. *et al.* Intracranial recordings from human auditory cortex reveal
1248 a neural population selective for musical song. *bioRxiv* (2019).
- 1249 55. Morosan, P. *et al.* Human primary auditory cortex: cytoarchitectonic subdivisions and
1250 mapping into a spatial reference system. *Neuroimage* **13**, 684–701 (2001).
- 1251 56. Baumann, S., Petkov, C. I. & Griffiths, T. D. A unified framework for the organization of
1252 the primate auditory cortex. *Frontiers in systems neuroscience* **7**, 11 (2013).
- 1253 57. de Cheveigné, A. & Parra, L. C. Joint decorrelation, a versatile tool for multichannel data
1254 analysis. *Neuroimage* **98**, 487–505 (2014).
- 1255 58. Murphy, K. P. *Machine learning: a probabilistic perspective*. (MIT press, 2012).
- 1256 59. Schoppe, O., Harper, N. S., Willmore, B. D., King, A. J. & Schnupp, J. W. Measuring the
1257 performance of neural models. *Frontiers in Computational Neuroscience* **10**, 10 (2016).
- 1258 60. Kell, A. J., Yamins, D. L., Shook, E. N., Norman-Haignere, S. V. & McDermott, J. H. A
1259 task-optimized neural network replicates human auditory behavior, predicts brain
1260 responses, and reveals a cortical processing hierarchy. *Neuron* (2018).
- 1261 61. Efron, B. & Tibshirani, R. J. *An introduction to the bootstrap*. (CRC press, 1994).
- 1262 62. Fisher, W. M. *tsylb: NIST syllabification software, version 2 revision 1.1*. (1997).
- 1263

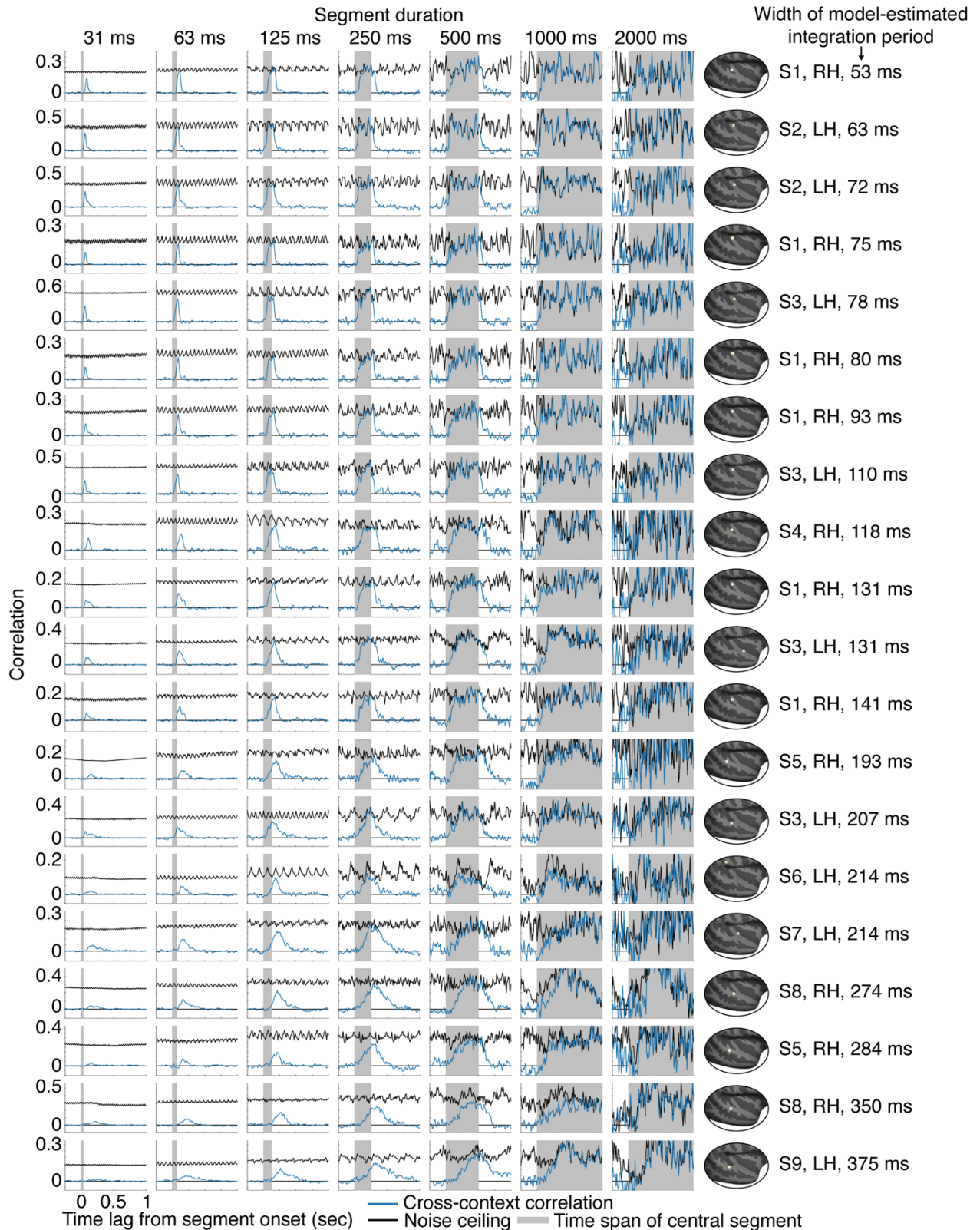
1264 **Acknowledgements**

1265 We thank Daniel Maksumov, Nikita Agrawal, Stephanie Montenegro, Leyao Yu, Marcin
1266 Leszczynski and Idan Tal for help with data collection. We thank Stephanie Montenegro and
1267 Hugh Wang for help in localizing electrodes. And we thank Alex Kell, Stephen David, Josh
1268 McDermott, Bevil Conway, Nancy Kanwisher, Nikolaus Kriegeskorte, and Marcin Leszczynski
1269 for comments on an earlier draft of this manuscript. This study was supported by the National
1270 Institutes of Health (NIDCD-DC014279, S10 OD018211, NINDS-R01-NS084142) and the
1271 Howard Hughes Medical Institute (LSRF postdoctoral award to SNH).



1272
1273
1274
1275
1276

Fig S1. Histogram of phoneme, syllable, and word durations in TIMIT. Durations of phonemes, multi-phoneme syllables, and multi-syllable words in the commonly used TIMIT database. Phonemes and words are labeled in the database. Syllables were computed from the phoneme labels using the software tsylb2⁶². The median duration for each structure is 64, 197, and 479 milliseconds, respectively.

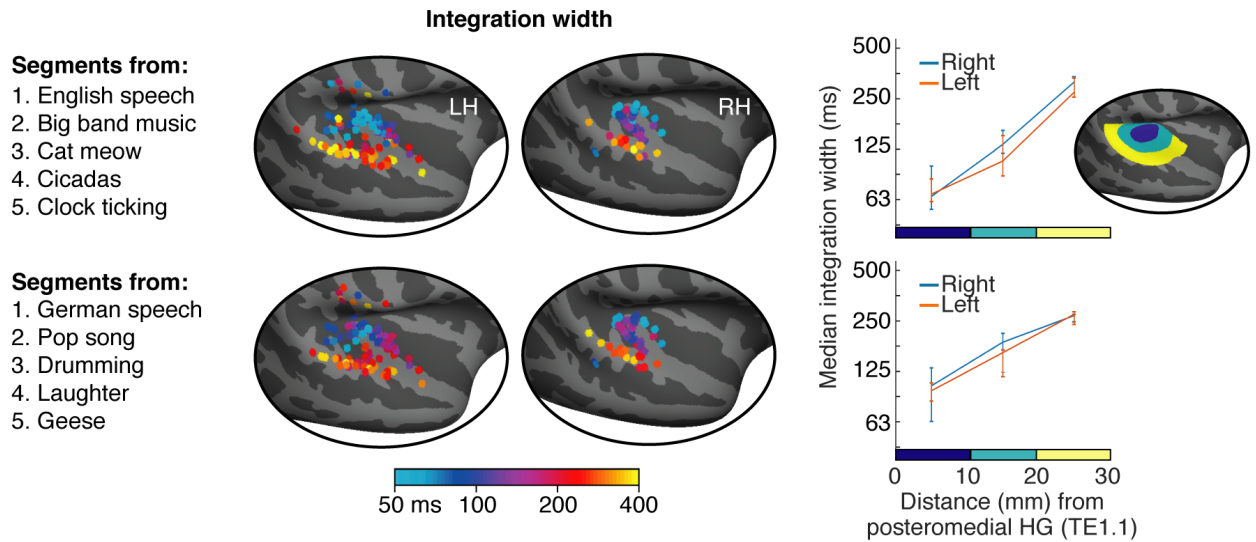


1277
1278
1279
1280
1281
1282
1283
1284
1285
1286
1287

Fig S2. Cross-context correlation for 20 representative electrodes. Electrodes were selected to illustrate the diversity of integration periods. Specifically, we partitioned all sound-responsive electrodes into 5 groups based on the width of their integration period, estimated using a model (Fig 3 illustrates the model). For each group, we plot the four electrodes with the highest SNR (as measured by the test-retest correlation across the sound set). Electrodes have been sorted by their integration width, which is indicated to the right of each plot, along with the location, hemisphere and subject index for each electrode. Each plot shows the cross-context correlation and noise ceiling for a single electrode and segment duration (indicated above each column). There were more segments for the shorter durations, and as a consequence, the cross-context correlation and noise ceiling were more stable/reliable for shorter segments (the number of segments is inversely proportional to the duration). This property is

1288 useful because at the short segment durations, there are a smaller number of relevant time lags, and
1289 it is useful if those lags are more reliable. The model used to estimate integration periods pooled across
1290 all lags and segment durations, taking into account the reliability of each datapoint (see *Model-*
1291 *estimated integration periods* in the Results and Methods).

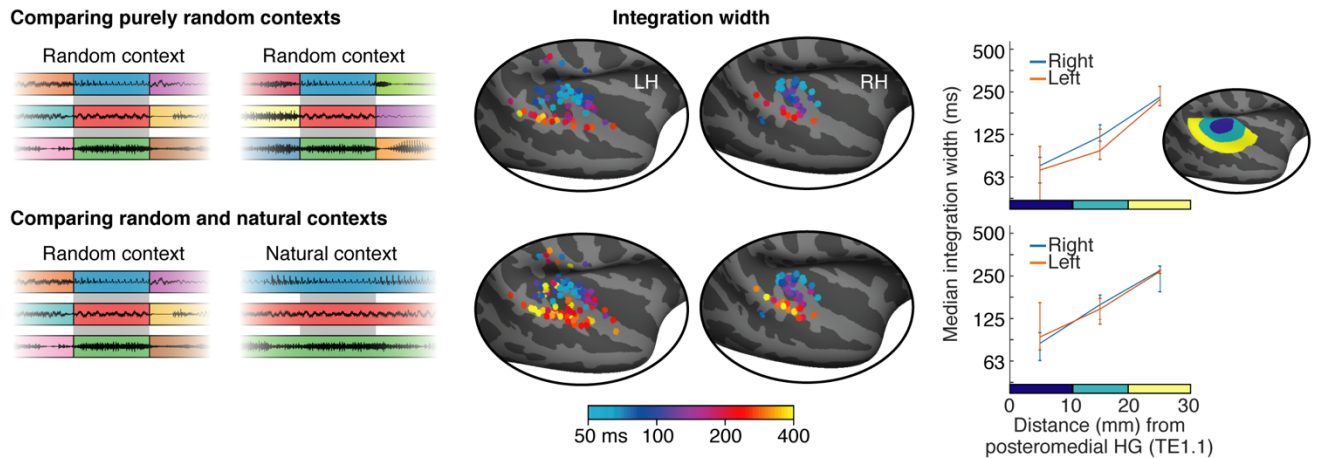
1292



1293
1294
1295
1296
1297
1298
1299

Fig S3. Split-half analysis across the sound set. Sound segments were excerpted from 10 sounds. To assess the robustness of our results to the sounds tested, we estimated integration periods using segments drawn from two non-overlapping splits of 5 sounds each (listed on the left). Since many non-primary regions only respond strongly to speech or music¹⁰⁻¹², we included speech and music in both splits. Format is analogous to **Figure 4** but only showing integration widths (integration centers were also similar between splits).

1300



1301

1302

1303

1304

1305

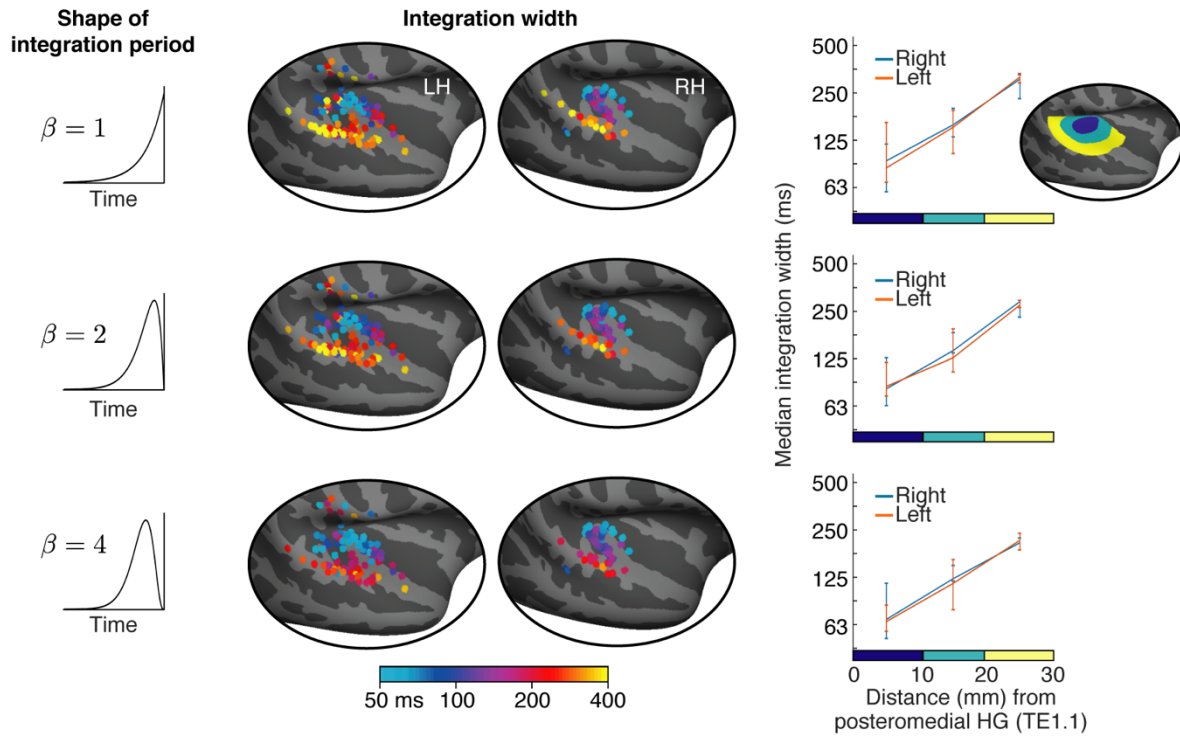
1306

1307

1308

1309

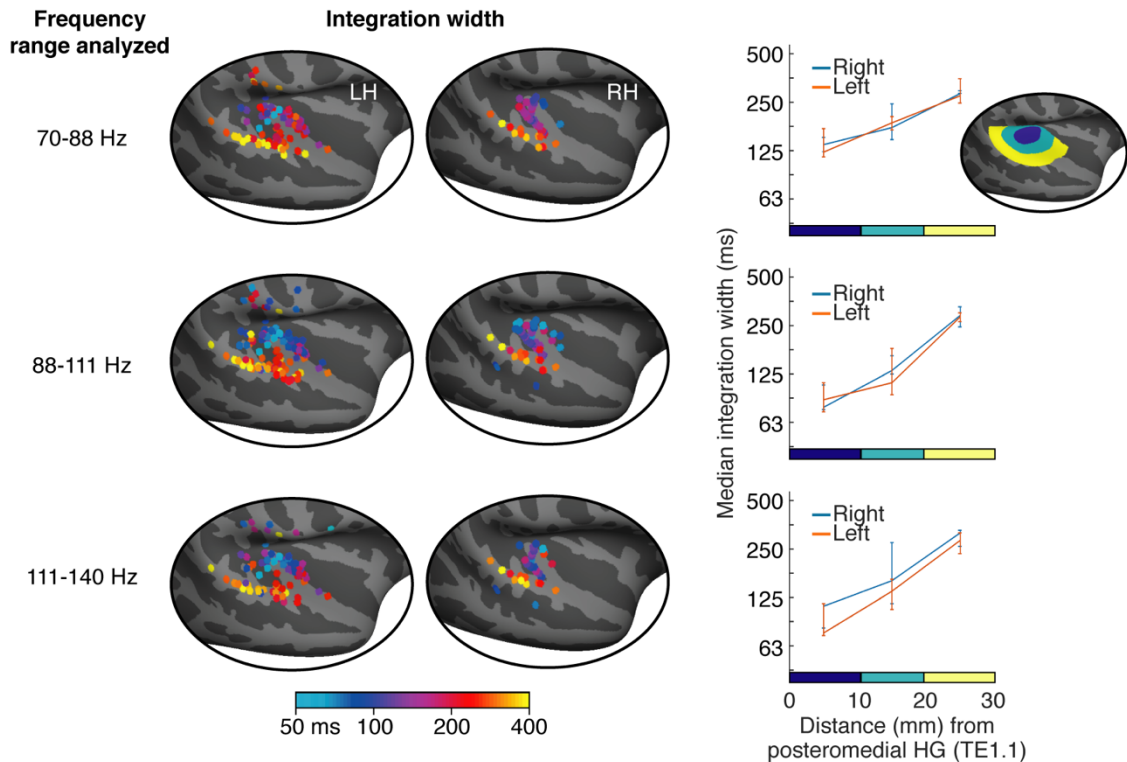
Fig S4. Comparing random and natural contexts. Shorter segments were created by subdividing longer segments, which made it possible to consider two types of context: (1) random context, in which each segment is surrounded by random other segments (2) natural context, where a segment is a subset of a longer segment and thus surrounded by its natural context. When comparing responses across contexts, one of the two contexts must always be random so that the contexts differ. But the other context can be random or natural. Our main analyses pooled across both types of comparison. Here, we show integration widths estimated by comparing either purely random contexts (top panel) or comparing random and natural contexts (bottom panel). Format is analogous to **Figure 4**.



1310
1311
1312
1313
1314
1315

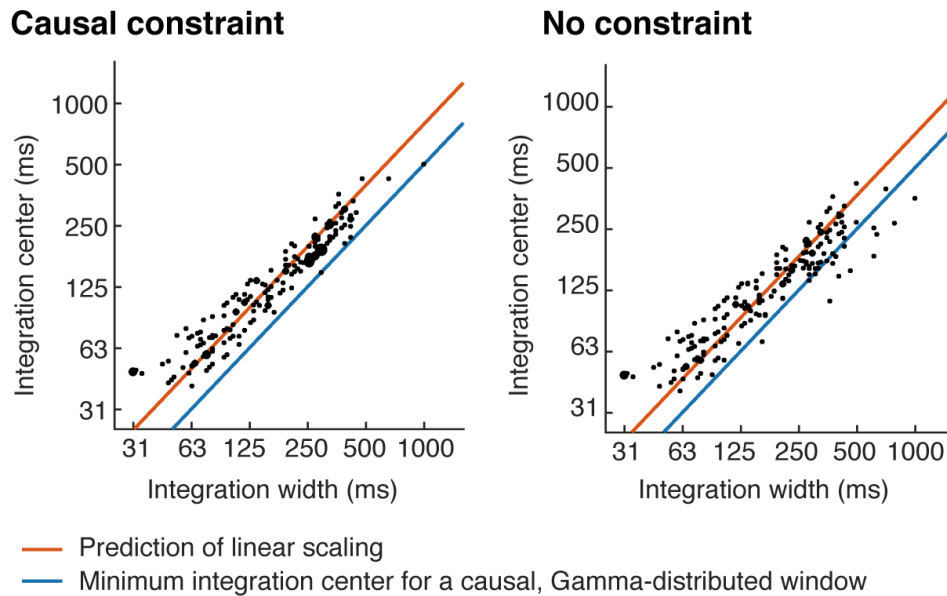
Fig S5. Results for different model window shapes. We modeled integration periods using window shapes that varied from more exponential to more Gaussian (the parameter β in equations 1&2 controls the shape of the window, see Methods). For our main analysis, we selected the shape that yielded the best prediction for each electrode. This figure plots integration widths estimated using three different fixed shapes. Format is analogous to **Figure 4**.

1316



1317
1318
1319
1320
1321
1322
1323
1324
1325

Fig S6. Results for different frequency ranges. For our primary analysis, we measured the broadband power of each electrode between 70 and 140 Hz. This figure shows the results of measuring integration widths from three different subsets of this broader range (equally spaced on a logarithmic scale). Format is analogous to **Figure 4**. Results were similar to those of our main analysis, but using a lower frequency range (70-88 Hz) appeared to limit the shortest integration widths that were detectable by our paradigm, plausibly because faster power fluctuations are better conveyed by a faster carrier.



1326

1327

1328

1329

1330

1331

1332

1333

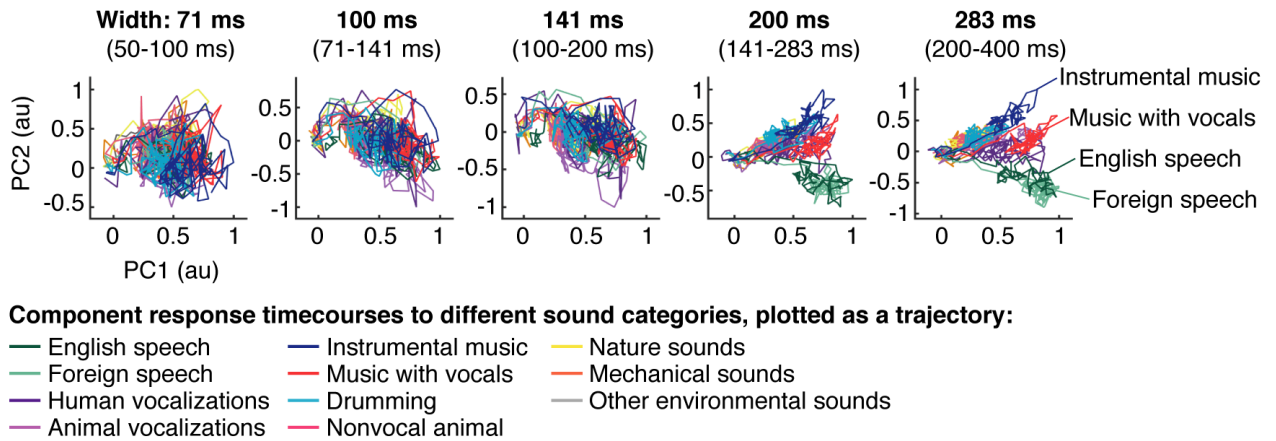
1334

1335

Fig S7. Relationship between integration widths and centers. This figure shows a scatter plot between integration centers and widths. Each dot corresponds to an electrode and larger dots indicate that multiple electrodes were assigned to that pairing of centers/widths. The integration width places a lower bound on the integration center for a causal window (blue line). Integration centers scaled approximately linearly with the integration width (orange line), and remained relatively close to the minimum possible for a casual window. On the left, we show results when integration periods were explicitly constrained to be causal, and on the right, we show results without this constraint. Results were similar in the two cases because the inferred integration periods were close-to-causal without being explicitly constrained to be so.

1336

Principle components at different integration periods



1337

1338

1339

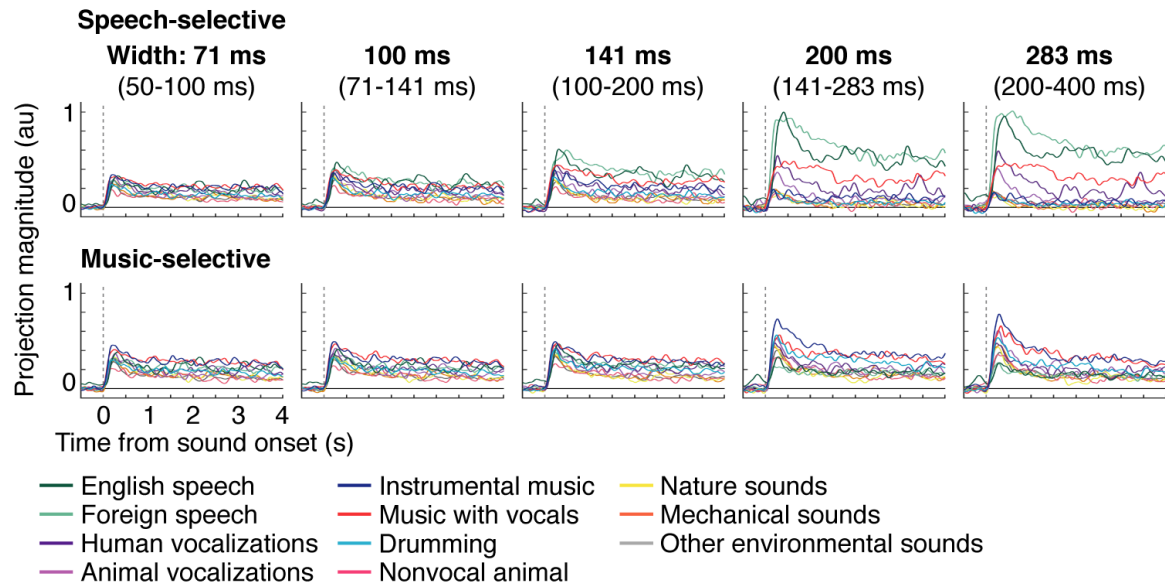
1340

1341

1342

Fig S8. Principal components at different integration periods. Electrodes were grouped based on the width of their integration period in octave intervals (shown above each plot). Responses were then projected onto the top two principal components from each group. This figure shows the average component response timecourse to each category, plotted as a trajectory. Format is the same **Figure 5a**.

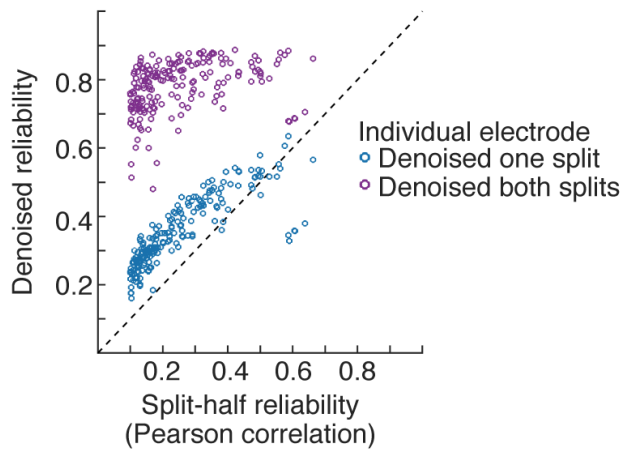
Response components most selective for speech and music



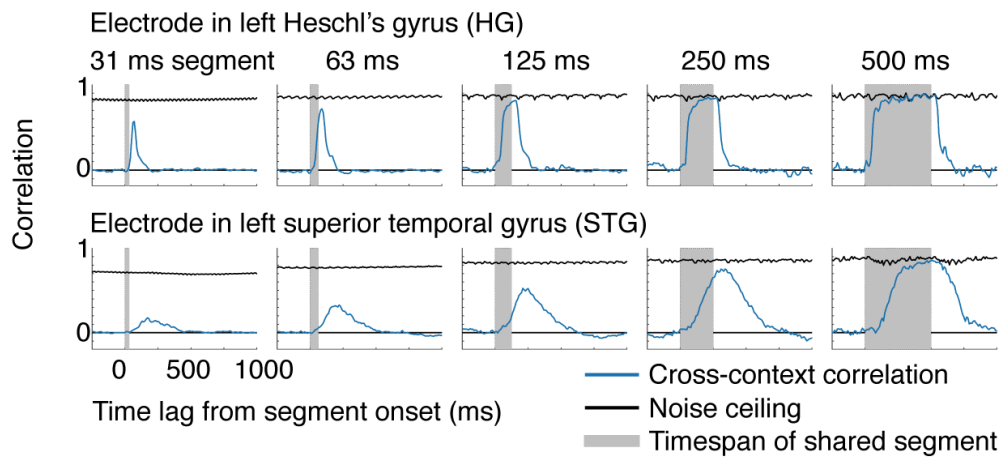
1343
1344 **Fig S9. Response components most selective for speech or music.** Electrodes were grouped
1345 based on the width of their integration period in octave intervals (shown above each plot). The
1346 responses from each group were then projected onto the components that showed the greatest speech
1347 (top panel) or music selectivity (bottom panel). The speech-selective component was optimized to
1348 separate responses to English and foreign speech from all other sounds (excluding vocal music which
1349 has speech). The music-selective component was optimized to separate responses to instrumental
1350 music, vocal music, and drumming from all other sounds. Each line reflects the average component
1351 response timecourse to one sound category. The timecourses have been smoothed so that lines for
1352 different categories can be clearly seen (100 millisecond FWHM Gaussian kernel). Independent sounds
1353 were used to estimate components and measure their response. **Figure 5c** quantifies how well
1354 separated speech and music are along each component. Separation was calculated using the
1355 response timecourses to individual sounds without any smoothing (using all timepoints between 250
1356 milliseconds and 4 seconds post-stimulus onset).

1357

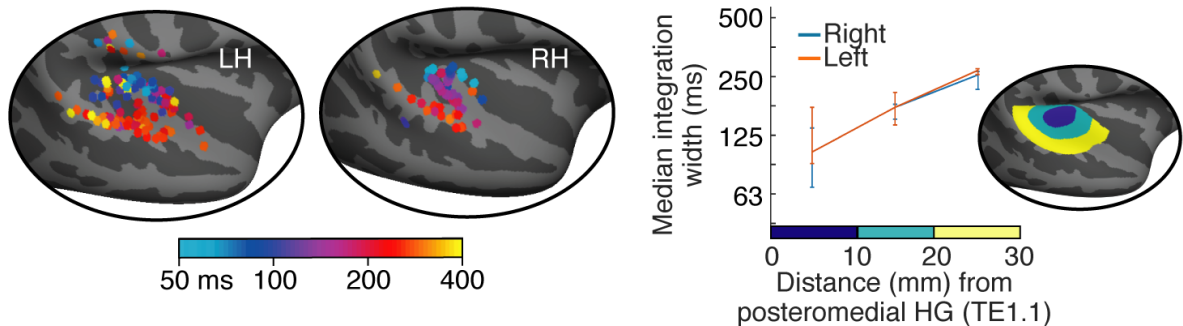
a Effect of denoising on data reliability



b Cross-context correlation computed from denoised responses



c Integration widths estimated from denoised responses



1358

1359

1360

1361

1362

1363

1364

1365

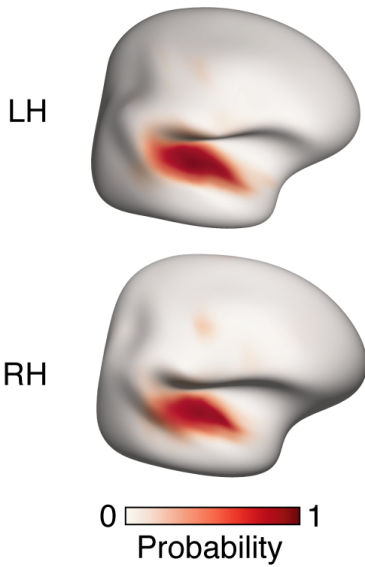
1366

1367

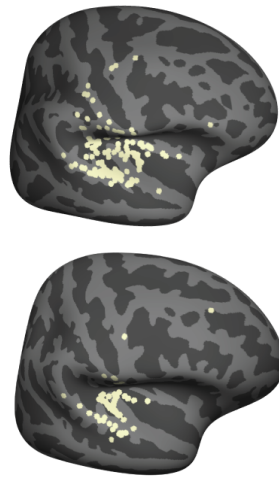
1368

Fig S10. Integration periods estimated from denoised data. **a**, Data were denoised by projecting electrode responses from one subject onto those from all other subjects. This procedure leads to a boost in SNR because the stimulus-driven response is more consistent across subjects than the noise. Each dot shows the split-half reliability of one electrode before (x-axis) or after (y-axis) denoising. The denoising procedure was either applied to both splits of data (purple dots) or to only one split of data (blue dots). Applying the analysis to both splits reveals the overall change in reliability. Applying the analysis to one split provides a fairer test of whether the denoising analysis removes more signal or noise. **b**, The cross-context correlation for the same example electrodes shown in **Fig 2b**, but measured from denoised responses. The trends are similar but the noise ceiling is much higher. **c**, Integration widths estimated from denoised responses. Same format as **Figure 4**.

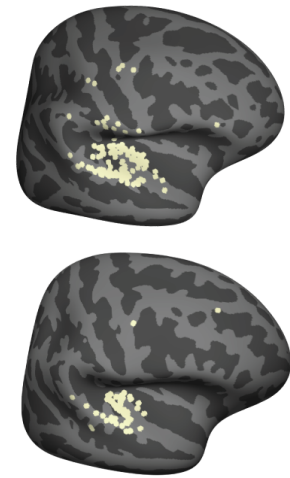
a Probability of a sound-driven response (measured with fMRI)



b Electrode localization without any constraint

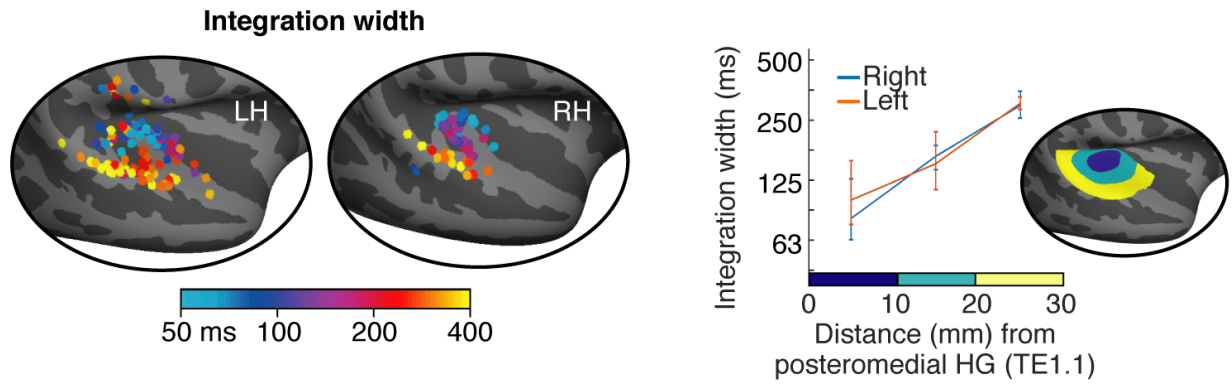


c Electrode localization constrained by probabilistic map of sound responses



1369
1370
1371
1372
1373
1374
1375
1376
1377
1378
1379
1380
1381

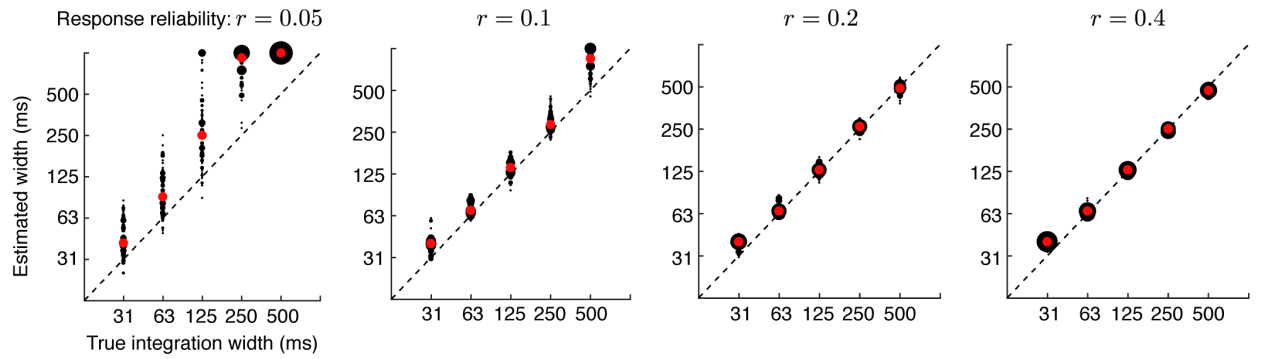
Fig S11. Constraining the anatomical localization of electrodes. **a**, Map showing the probability of observing a significant response to sound at each point in the brain. The map was computed using whole-brain fMRI responses to large collection of 192 natural sounds across a large cohort of 20 subjects³⁹. **b**, Electrode localization based on mapping each electrode to the nearest point on the cortical surface. Due to cortical folding, nearby points in space can be faraway on the cortical surface. As a consequence, small localization errors can cause electrodes to be mapped to the wrong region. Such errors like explain why some electrodes have been localized to the supramarginal gyrus, which abuts the superior temporal gyrus where responses to sound are much more common. **c**, To minimize gross localization errors, we treated the probability map of sound-driven responses shown in panel A as a prior and used to it constrain the localization (see *Electrode localization* in the Methods). Because the prior map is highly smooth this approach did not substantially affect the localization of electrodes at a fine spatial scale.



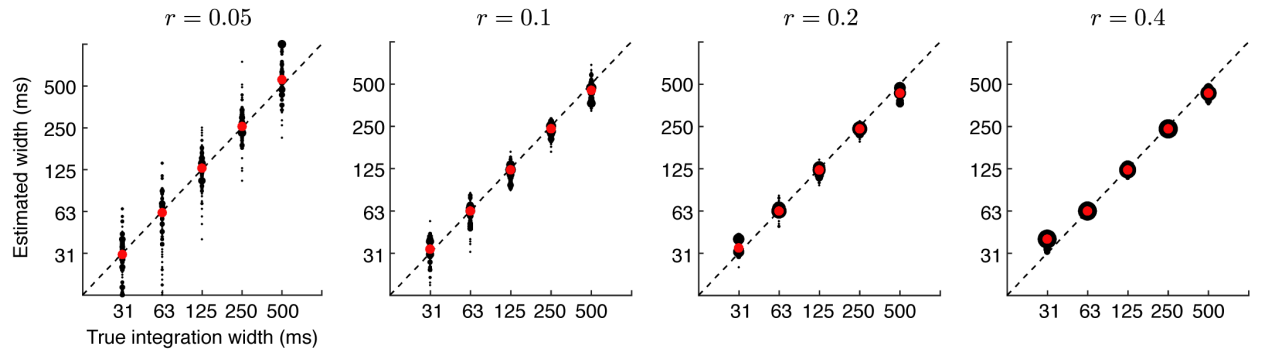
1382
1383
1384
1385
1386
1387

Fig S12. Results using the squared error uncorrected for bias. For our main analysis we quantified prediction accuracy using a bias-corrected variant of the squared error. Here, we plot integration widths estimated using the uncorrected squared error. Results were similar using the uncorrected and bias-corrected error (compare with **Fig 4**), suggesting our data were sufficiently reliable that the bias had little effect.

Squared error metric



Bias corrected



1388

1389

1390 **Fig S13. Results of model simulation.** We simulated a response that integrated sound amplitude
1391 within a Gamma-distributed integration period. The integrated amplitudes were then used to modulate
1392 the power of a broadband gamma signal. We tested our ability to infer the correct integration width
1393 using our complete analysis pipeline. Black dots show the estimated width for a single simulated
1394 response, and red dots show the median width across 100 simulations. Results are plotted for different
1395 SNRs, manipulated by adding variable amounts of noise to achieve a desired test-retest reliability (the
1396 split-half correlation of the simulated data is shown above each plot). When using the squared error to
1397 measure prediction accuracy, we found there was an upward bias for low SNRs (top panel). To address
1398 this issue, we derived a variant of the squared error that largely corrected the bias (bottom panel) (see
Bias correction in Methods).

1399
1400

Table S1. Sound sources used for TCI experiment.

English speech
German speech
Big band music
Pop song
Drum solo
Laughter
Cat meows
Geese
Cicadas
Clock ticking

1401