
Research Article: New Research | Cognition and Behavior

Distinguishing fine structure and summary representation of sound textures from neural activity

<https://doi.org/10.1523/ENEURO.0026-23.2023>

Cite as: eNeuro 2023; 10.1523/ENEURO.0026-23.2023

Received: 23 January 2023

Revised: 25 August 2023

Accepted: 31 August 2023

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

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

Copyright © 2023 Berto et al.

This is an open-access article distributed under the terms of the Creative Commons Attribution 4.0 International license, which permits unrestricted use, distribution and reproduction in any medium provided that the original work is properly attributed.

Distinguishing fine structure and summary representation of sound textures from neural activity

Abbreviated title: Neural response to auditory details and statistics

Author names

Martina Berto¹, Emiliano Ricciardi¹, Pietro Pietrini¹, Nathan Weisz^{2,3}, Davide Bottari¹

Affiliations

- 1 Molecular Mind Lab, IMT School for Advanced Studies Lucca, Lucca, Italy
- 2 Department of Psychology and Centre for Cognitive Neuroscience, University of Salzburg, Austria
- 3 Neuroscience Institute, Christian Doppler University Hospital, Paracelsus Medical University, Salzburg, Austria

Author Contributions

MB and DB Designed research; MB Performed research; MB and DB Analyzed data; MB, DB, ER, PP, and NW Wrote the paper.

Corresponding author

Martina Berto; martina.berto@imtlucca.it

Number of figures: 3

Number of tables: 0

Number of Multimedia: 3

Number of words for Abstract: 213

Number of words for Significance Statement: 98

Number of words for Introduction: 971

Number of words for Discussion: 2318

Acknowledgments

The authors thank all the students who helped with recruiting participants and/or data collection: Nicolò Castellani, Irene Sanchez, Chiara Battaglini, and Dila Suay.

Conflict of interest

Authors declare no conflict of Interest

Funding sources:

Davide Bottari (PRIN 2017 research grant. Prot. 20177894ZH).

1 **ABSTRACT**

2 The auditory system relies on both local and summary representations; acoustic local
3 features exceeding system constraints are compacted into a set of summary statistics.
4 Such compression is pivotal for sound-object recognition. Here, we assessed whether
5 computations subtending local and statistical representations of sounds could be
6 distinguished at the neural level. A computational auditory model was employed to
7 extract auditory statistics from natural sound textures (i.e., fire, rain) and to generate
8 synthetic exemplars where local and statistical properties were controlled. Twenty-four
9 human participants were passively exposed to auditory streams while the EEG was
10 recorded. Each stream could consist of short, medium, or long sounds to vary the
11 amount of acoustic information. Short and long sounds were expected to engage local
12 or summary statistics representations, respectively. Data revealed a clear dissociation.
13 Compared to summary-based ones, auditory-evoked responses based on local
14 information were selectively greater in magnitude in short sounds. Opposite patterns
15 emerged for longer sounds. Neural oscillations revealed that local features and
16 summary statistics rely on neural activity occurring at different temporal scales, faster
17 (beta) or slower (theta-alpha). These dissociations emerged automatically without
18 explicit engagement in a discrimination task. Overall, this study demonstrates that the
19 auditory system developed distinct coding mechanisms to discriminate changes in the
20 acoustic environment based on fine structure and summary representations.

21

22

23

24 **SIGNIFICANCE STATEMENT**

25 Prior to this study, it was unknown whether we could measure auditory discrimination
26 based on local temporal features or spectrotemporal statistics properties of sounds from
27 brain responses. Results show that the two auditory modes of sound discrimination
28 (local and summary statistics) are automatically attuned to the temporal resolution (high
29 or low) at which a change has occurred. In line with the temporal resolutions of auditory
30 statistics, faster or slower neural oscillations (temporal scales) code sound changes
31 based on local or summary representations. These findings expand our knowledge of
32 some fundamental mechanisms underlying the function of the auditory system.

33

34

35 **INTRODUCTION**

36 The human auditory system can discriminate sounds at both high and low temporal
37 resolutions (McAdams, 1993; Griffiths, 2001). The processing of fine temporal details
38 relies on extracting and retaining local acoustic features (on the order of a few
39 milliseconds) to detect transient changes over time (Plomp, 1964; McDermott,
40 Schemitsch, and Simoncelli, 2013; Dau, Kollmeier, and Kohlrausch, 1997). These
41 temporal variations characterize different sound objects and help the system discern
42 among acoustic sources. However, environmental inputs typically comprise long-lasting
43 sounds in which the number of local features to be retained exceeds the sensory
44 storage capacity. For this reason, the system may need to condense information into
45 more compact representations to discriminate sounds over more extended periods

46 (McDermott, Schemitsch, and Simoncelli, 2013). As the duration of the entering sounds
47 increases, summary representations are built upon fine-grained acoustic features to
48 condense information into a more compact and retainable structure (Yabe et al., 1998).
49 The processing of summary representations allows abstraction from local acoustic
50 features and prompt sound categorization (McDermott and Simoncelli, 2011;
51 McDermott, Schemitsch, and Simoncelli, 2013).

52 For sounds characterized by a constant repetition of similar events over time (such as
53 sound textures, e.g., rain, fire, typewriting; Saint-Arnaud and Popat, 1995), this form of
54 compression consists of a set of auditory statistics comprising averages over time of
55 acoustic amplitude modulations at different frequencies (McDermott and Simoncelli,
56 2011; Figure 1A).

57 Computational approaches in auditory neuroscience allow the mathematical
58 formalization of this set of auditory statistics (Figure 1A). The basic assumption is
59 derived from information theories (Barlow, 1961) and suggests that if the brain
60 represents sensory input with a set of measurements (statistics), any signal containing
61 values matching those measurements will be perceived as the same.

62 Psychophysical experiments revealed that stimuli including the same summary statistics
63 -but different local features- are easy to discriminate when they are short, but that as
64 duration increases and summary representation takes over, they are progressively more
65 challenging to tell apart (Berto et al., 2021; McDermott, Schemitsch, and Simoncelli,
66 2013). On the other hand, when sounds comprise different statistics, their perceived
67 dissimilarity will increase with duration as their summary representations diverge (Berto
68 et al., 2021; McDermott, Schemitsch, and Simoncelli, 2013). While some evidence

69 exists in the animal model (see Zhai et al., 2020, for results in rabbits), the neural
70 activity underpinning local features and summary statistics is unknown in humans.
71 Moreover, previous behavioral studies required participants to attend to stimuli to
72 perform a task actively. From this evidence alone, it thus remains unanswered whether
73 discrimination based on local features and their summary statistics can occur despite
74 the lack of an active task and can therefore occur automatically.

75 To fill these gaps, we used a validated computational auditory model (McDermott and
76 Simoncelli, 2011) to extract auditory summary statistics from natural sounds and
77 generate synthetic sounds that feature this same set of measurements (see Material
78 and Methods; Figure 1A,B). With this approach, it is possible to impose the same set of
79 statistics on different white noise samples that initially had different local structures
80 (Figure 1B,C). By employing this synthesis approach, we could create sounds that
81 differ at high temporal resolutions (e.g., local features) but are perceptually
82 indistinguishable at lower ones (summary statistics) and vice versa (Figure 1C). We
83 acquired EEG measurements in participants passively exposed to streams composed of
84 triplets of sounds presented at a fast stimulation rate (2Hz). To ensure generalizability,
85 sounds were randomly drawn from a large set of synthetic excerpts (see Material and
86 Methods). Within each triplet, the first sound was repeated twice, while the third one
87 was novel. Two experiments were designed (Figure 2A). (1) In Local Features, the
88 novel and repeated sounds differed only in their local structures, as they were
89 generated by imposing the same auditory statistics on different white noise samples; (2)
90 in Summary Statistics, the novel sound was generated from the same white noise
91 sample but differed from the repeated ones as it comprised a different set of auditory

92 statistics. As summary statistics are expected to be relevant at increasing sound
93 duration (McDermott, Schemitsch, and Simoncelli, 2013), sounds including the same
94 statistics but originating from different input white noises will be easily distinguishable at
95 short duration but not at long ones (Figure 1D). By contrast, sounds derived from the
96 same white noise sample but including different summary statistics will have different
97 statistical values when measured at long durations but more similar values when
98 measured at short durations (Figure 1D). In fact, at short durations, statistics will be
99 influenced by their similar temporal structure (see Figure 1F).

100 Thus, to manipulate the extent of temporal and statistical similarity, we presented
101 separate sound streams comprising stimuli of different lengths (either 40, 209, or
102 478ms; Figure 2A). First, we investigated auditory-evoked responses to uncover
103 magnitude changes in neural activity associated with the two modes of representation.
104 We predicted that short and long sounds would prompt larger auditory-discriminative
105 responses for local features and summary statistics, respectively. Specifically, we
106 hypothesized that since the amount of information (e.g., sound duration) impacts the
107 statistical similarity of sound excerpts, distinct mechanisms are engaged in the
108 processing of local features compared to summary statistics emerging over time. That
109 is, in the case of short sounds, the brain may emphasize transient amplitude
110 modulations (i.e., broadband envelope changes), while spectrotemporal statistics will
111 become informative as sound size increases.

112 In line with this prediction, we expected brief local information to be encoded at a faster
113 timescale (Panzeri et al., 2010) than summary statistics. That is, we expect the
114 response pattern of the neuronal populations involved in processing local features to be

115 encoded at higher frequency ranges and earlier latencies of neural oscillations
116 compared with summary statistics. To this end, we investigated neural oscillations and
117 assessed whether information measured at different temporal scales in the oscillatory
118 pattern revealed specific fingerprints of discrimination based on local features and
119 summary statistics.

120

121 **MATERIALS AND METHODS**

122 **Participants**

123 Twenty-four normal-hearing right-handed young adults (12 of either sex; mean age=
124 27.13 years, std= 2.83) participated in the experiment. All participants were healthy;
125 they were fully informed of the scope of the experiment, signed written informed consent
126 before testing, and received monetary compensation. The study was approved by the
127 regional ethical committee, and the protocol adhered to the guidelines of the Declaration
128 of Helsinki (2013).

129

130 **Sample size estimation**

131 This sample size was estimated via simulations. We used the procedure described in
132 Wang and Zhang (2021) and simulated a dataset with two conditions (Local Features
133 and Summary Statistics) of Auditory Evoked Potentials data. First, we selected three
134 electrodes of interest at central locations (E7, E65, E54). For the simulation, we chose a
135 time window between 0.1 and 0.3s based on previous MMN studies (see Näätänen et
136 al., 2007 for review). The amplitude values at the electrodes of interest for the two
137 conditions were sampled from a bivariate normal distribution (within-subject design)

138 where mean and standard deviation were chosen based on results of four pilot datasets
139 (mean Local Features= 0.16; mean Summary Statistics= 0.56; std Local Features=
140 0.52; std Summary Statistics= 0.54).

141 We then ran a cluster-based permutation on simulated datasets to test whether any
142 statistical cluster (t-values) exhibited a significant difference between the two conditions
143 with an alpha level of 0.05. The procedure started with a sample size of 10 and
144 increased in steps of one until it reached a power of 0.80. We ran 1000 simulations for
145 each sample size and calculated the power as the proportion of the number of times
146 significant clusters were found in these 1000 simulations. The simulation results showed
147 that to obtain power above 0.8, a sample size of $N= 24$ was required.

148 The algorithm to perform such analyses can be downloaded from this link:

149 <https://osf.io/rmqhc/>

150

151 **Stimuli**

152 Synthetic sounds were generated using a previously validated computational auditory
153 model of the periphery. The auditory model and synthesis toolbox are available at:
154 <http://mcdermottlab.mit.edu/downloads.html>.

155 This auditory model emulates basic computations occurring in the cochlea and midbrain
156 (McDermott and Simoncelli, 2011).

157 The signal (7s original recording of a sound texture, $N=54$; see Extended Data Table 1-
158 2) was decomposed into 32 cochlear subbands using a set of gammatone filter banks
159 with different central frequencies spaced on an ERB scale. Absolute values of the
160 Hilbert transform for each subband were computed to extract the envelope modulation

161 of each cochlear channel over time. Envelopes were then compressed to account for
162 the nonlinear transformations performed by the cochlea. The first set of statistics was
163 measured from the transformed envelopes: mean, skewness, variance, autocorrelation
164 (within each cochlear channel), and cross-correlation (between channels). Additional
165 filtering was applied to the envelopes to account for the modulatory response of the
166 spectrotemporal receptive fields of neurons in the midbrain (Bacon and Wesley
167 Grantham, 1989; Dau et al., 1997). Three additional statistics resulting from these
168 operations could be derived: modulation power, C1, and C2 (respectively, the
169 correlation between different envelopes filtered through the same modulation filter and
170 the correlation between the same envelopes filtered by other modulation filters; Figure
171 1A). The resulting set of statistics extracted from the original recording of sound textures
172 was imposed on four 5s white noise samples (Figure 1A, B, C). This allowed the
173 generation of four different sound exemplars for each sound texture, which varied
174 selectively in their local features but included similar long-term summary
175 representations (Figure 1C). All synthetic exemplars featuring the same auditory
176 statistics were perceptually very similar to the original sound texture from which they
177 were derived, even when their input sounds (white noise) varied (Figure 1C-E).
178 Synthetic sounds with the same imposed auditory statistics represent different
179 exemplars of the same sound texture with the same summary statistics but a different
180 fine-grained structure. This is because, in the synthesis procedure, the imposed
181 statistics are combined with the fine structure of the original white noise sample (Figure
182 1B).

183 Importantly, to create experimental stimuli, all four 5s synthetic exemplars were cut from
184 the beginning to the end into excerpts of different lengths, either short (40ms), medium
185 (209ms) or long (478ms). These lengths were chosen based on results in previous
186 behavioral investigations (Berto et al., 2021; McDermott, Schemitsch, and Simoncelli,
187 2013). Excerpts were equalized to the same root mean square amplitude (RMS= 0.1)
188 and had a sampling rate of 20kHz. A 20ms ramp (half-hann window) was applied to
189 each excerpt, 10ms at the beginning and 10ms at the end, to avoid edge artifacts
190 (McDermott, Schemitsch, and Simoncelli, 2013). The stimuli used here were validated
191 in a previous study (Berto et al., 2021) in which we replicated the original finding
192 (McDermott, Schemitsch, and Simoncelli, 2013). The experimental stimuli presented for
193 each run were randomly drawn from all available excerpts according to the experiment
194 requests (see below).

195

196 **Procedure**

197 Participants were tested in a sound-isolation booth. After reading instructions on a
198 monitor, they listened to the sounds in the absence of retinal input (participants were
199 blindfolded to prevent visual input).

200 For each run of the experimental session, a sound sequence lasting 108s was
201 presented. The series contained triplets of sounds ($n = 216$) presented one after the
202 other to form an almost continuous sound stream, in which sound onsets occurred
203 every 500ms (Figure 2A). Within each sequence, all sounds had the same duration
204 (either 40, 209, or 478ms).

205 Two experiments were implemented: (1) In Local Features, two different 5s synthetic
206 exemplars of the same sound texture were selected (out of the four we had created);
207 the combination of selected pair of exemplars vary randomly across triplets (e.g., first
208 and second; second and fourth, and so on). These two exemplars were cut into brief
209 excerpts of either 40, 209, or 478ms. According to the selected duration (which was
210 different for each sequence), two excerpts (one for each exemplar) were chosen from
211 among the available ones. The two excerpts had the same starting point (in seconds)
212 from the onset of the 5s exemplar. The first sound excerpt was repeated twice, and
213 afterward, the other was presented as the third element in the triplet.

214 Thus, two sounds within a triplet were identical (repeated), while the third one (novel)
215 comprised different local features but converging summary statistics; in other words,
216 repeated and novel sounds had the same generative statistics (both could be, e.g.,
217 waterfall) but different acoustic local features (Figure 2A, left panel; Extended Data
218 Table 1-2, column 1). (2) In Summary Statistics, sound textures were coupled according
219 to their perceived similarity (McDermott, Schemitsch, and Simoncelli, 2013; see
220 Extended Data Table 1-2, columns 1 and 2). For the textures in column 1, one out of the
221 four 5s synthetic exemplars was selected and cut into excerpts of the required duration
222 (40, 209, or 478ms); one of such excerpts was picked randomly. The same was done
223 for the coupled texture, ensuring that both were derived from the same white noise
224 sample and that both drawn excerpts had the same starting point in seconds. Thus, we
225 ensured that the sounds came from the same segment of the original input noise
226 sample and varied only for their imposed statistics. Again, the first excerpt was repeated
227 twice, while the other was used as the last sound in the triplet. The novel sound thus

228 deviated from the other two (repeated) in its auditory statistics, as it was a segment of
229 an exemplar of a different sound texture. This means the novel sound was a different
230 sound object (e.g., the repeated sounds might be waterfall excerpts and the novel one
231 air conditioner; see Figure 2A, right panel). However, since both originated from the
232 same segment of the same input white noise sample, their temporal structure (i.e.,
233 broadband envelope) measured at high resolution (that is, in brief excerpts) was
234 expected to be more similar in Summary Statistics than in the Local Features
235 experiment. This was indeed the case (see Figure 1F) and would affect the similarity of
236 statistics measured from short (but not long) sound excerpts (Figure 1D).

237 To ensure generalizability, the sound textures were different across triplets, so the
238 statistical similarity between repeated and novel sounds was kept constant within an
239 experiment while presenting different types of stationary sound objects.

240 Discriminative responses emerging from the contrast between the novel and repeated
241 sounds did not depend on specific properties (e.g., a change in frequency between a
242 particular type of sound category) but only on their local or statistical changes.

243 In both experiments, the order of the triplets was shuffled for each participant and run.
244 Moreover, excerpts were selected randomly from among those that shared the required
245 characteristics, so not only the presentation order but also stimuli per se were always
246 different across participants.

247 A total of six conditions were employed: two experiments (Local Features and Summary
248 Statistics) for three sound durations (40, 209, 478ms). Note that for each sound texture,
249 we synthesized only four exemplars that we cut into excerpts of different sound
250 durations (short, medium, or long). This means that within one experiment, the

251 presented excerpts belonged to the same pool of synthetic sounds, and only their
252 duration changed, not their properties. Thus, any dissociation between experiments
253 (Local Features or Summary Statistics) according to sound duration would indicate that
254 the processing of either local features or summary statistics strictly depends on the
255 amount of information presented.

256 Two sequences/runs per condition (Experiment * Duration) were presented for a total of
257 twelve runs. The order of runs was randomized across participants, and short breaks
258 were taken between runs. In a sound stream, excerpts were presented in triplets, with
259 the repeated one presented twice, followed by the novel one. Keeping the number of
260 repeated sounds constant allowed to control for the effects that differences in their
261 number could have on the brain response (e.g., standard formation, the effect by which
262 the number of repeated stimuli influence the response to the deviant element; see
263 Sussman and Gumejuk, 2005); moreover, it allowed to keep the duration of the
264 streams constant while manipulating the amount of information they encompass (e.g.,
265 the size of each sound excerpts). On the other hand, by keeping the novel position fixed
266 (as the third element of the triplet), we controlled for between-experiment differences in
267 expectancy effects (e.g., some novel sounds could be more predictable than others at
268 specific durations or based on their intrinsic properties) and, more importantly, we
269 ensured that the novel sound varied from the repeated ones only for its generative
270 statistics (same in Local Features and different in Summary Statistics) and original fine
271 structure (different in Local Features and same in Summary Statistics).

272 Since the interstimulus gap always depended on sound duration (sound onset was kept
273 constant at every 500ms), comparisons were assessed between experiments and within
274 the duration.

275 Participants had to listen to the sound stream but were asked to perform an orthogonal
276 task consisting of pressing a button when a beep sound was heard. The beep was a
277 pure tone higher in pitch and intensity than the sound-texture stream. The pure tone
278 was 50ms in length, had a frequency of 2200Hz, an amplitude of 50dB, a sampling rate
279 of 20kHz, and an RMS of 5. The beeps randomly occurred during the stimulation period.
280 The number of beeps varied randomly across runs from 0 to 3. Detection was
281 considered valid when the participant pressed the key within an arbitrary window of 3s
282 from beep occurrence.

283

284 **Similarity of summary statistics as a function of sound duration**

285 In order to assess the impact of sound duration on the statistical similarity between pairs
286 of excerpts, we extracted statistical values from all the pairs of excerpts (repeated and
287 novel) presented in the experiment to all participants and in all runs ($n= 20736$; note that
288 stimuli would appear more than once, as we adhered to the exact sound sequences
289 presented to participants). That is, for each synthetic excerpt pair, we extracted the set
290 of summary statistics (envelope mean, skewness, variance, and cross-band correlation;
291 modulation power, C1, and C2) through the auditory texture model (Figure 1A;
292 McDermott and Simoncelli, 2011). To assess similarity between summary statistic of
293 repeated and novel sounds, we used a similar procedure to the one employed during
294 sound synthesis to evaluates the quality of the output. This procedure consists of

295 computing the signal-to-noise ratio (SNR) between statistic classes measured from the
 296 synthetic signal and the original sound texture (McDermott and Simoncelli, 2011).
 297 Firstly, we computed the total squared error ε of statistics measured from repeated
 298 sounds and the corresponding novel sound at each cochlear channel k ($n=32$) as
 299 follow:

$$\varepsilon_k = (StatRep - StatNov)^2 ,$$

$$k \in [1,2,3 \dots, 32]$$

300 where *StatRep* is a statistic class (i.e., envelope mean, variance, or modulation power)
 301 measured from a repeated sound excerpt and *StatNov* is the same statistic class
 302 measured from the corresponding novel sound in the triplet. Note that for statistic
 303 classes that had more than the one dimension k (i.e., modulation power and
 304 correlations) the values across other dimensions (i.e., modulation bands) were summed
 305 prior to compute the error, as in McDermott and Simoncelli (2011).
 306 Secondly, we calculated the SNR for each statistic class by dividing the sum of the
 307 squared statistic values measured from the repeated sound by the squared error
 308 between repeated and novel sounds as follow:

$$SNR = 10 \log_{10} \left(\frac{\sum_k StatRep(k)^2}{\sum_k \varepsilon(k)} \right), k \in [1,2,3 \dots, 32]$$

309 We computed one SNR for each statistic class ($n=7$) and then average their values to
 310 have one average SNR for each excerpt pair presented in each experiment and
 311 duration. Average SNRs are displayed in Figure 1D.
 312 We then compared whether the average SNRs of sound excerpts were significantly
 313 different between experiment and within duration by performing non-parametric tests
 314 (Wilcoxon rank sum test). The results showed a clear dissociation according to sound

315 duration. When sounds were short (40ms), the average SNR of statistics between
316 repeated and novel sounds was higher in the Summary Statistics experiment ($p <$
317 0.001 , mean=9.94; std=2.4) than in the Local Features one (mean= 8.34; std= 1.24).
318 Namely, when sounds were short, statistical values were influenced by the white noise
319 sample, thus sounds originated from the same seed had more similar values compared
320 to when they originated from a different one, disregarding the generative statistics that
321 were imposed. Thus, we expected larger neural discriminatory responses in Local
322 Features experiment compared to the Summary Statistics one.

323 Conversely, at long duration (478ms), the average statistic SNR between repeated and
324 novel sounds was more dissimilar in the Summary Statistics experiment ($p < 0.001$,
325 mean= 6.73; std= 2.10) than in Local Features one (mean=9.31; std=1.23). At
326 increasing sound duration, summary statistics were no longer influenced by the
327 temporal structure of the original white noise sample as they converged to their original
328 values. Based on this observation, we expected greater neural activation in response to
329 Summary Statistics change compared to Local Features when sounds were long. The
330 same pattern was observed for medium sound duration (209ms; $p < 0.001$, mean
331 Summary Statistics= 8.00; std= 2.21; mean Local Features=9.19; std=1.30), although
332 there was a clear trend of decreasing average SNR with increasing sound duration in
333 the Summary Statistics experiment (see Figure 1D).

334 Overall, this analysis showed that the statistical similarity measured from the presented
335 sounds well predicted the brain response observed in the EEG.

336

337 **Similarity of temporal amplitude modulation in brief excerpts**

338 The previous analysis showed higher statistical similarity measured at high (but not low)
339 temporal resolutions from the excerpt pairs presented in the Summary Statistics
340 experiment. To test the hypothesis that this effect depended on the original temporal
341 structure of white noise samples (which will be more similar in the Summary Statistics
342 experiment compared to the Local Features one), we conducted a similar correlation
343 analysis for brief excerpts, but this time using as dependent variables the excerpts
344 broadband amplitude modulations and disregarding their spectral density. Specifically,
345 for every sound pair presented across participants, we used the auditory texture model
346 (Figure 1A; McDermott and Simoncelli, 2011) to compute the cochleograms of all the
347 40ms excerpts presented in the study ($n = 6912$) and averaged them across frequency
348 bands to extract their broadband envelopes. We then computed Pearson's correlations
349 between the envelopes of each excerpt pair (repeated and novel) to estimate their linear
350 relationship (Figure 1F). The correlation coefficients (r) were transformed into Fisher- z
351 scores for statistical comparison by t -tests. The results showed that the amplitude
352 modulations over time between excerpt pairs were more correlated in the Summary
353 Statistics experiment (mean = 2.35, std = 0.71) than in the Local Features experiment
354 (mean = 1.63, std = 0.51, p -value < 0.001). This result confirmed that regardless of their
355 spectral density, the repeated and novel sounds in the Summary Statistics experiment
356 shared more comparable temporal amplitude modulations than those in the Local
357 Features experiment.

358

359 **EEG recording**

360 Electroencephalography (EEG) was recorded from an EGI HydroCel Geodesic Sensor
361 Net with 65 EEG channels and a Net Amps 400 amplifier (Electrical Geodesics, Inc.,
362 EGI, USA). The acquisition was obtained via EGI's Net Station 5 software (Electrical
363 Geodesics, Inc., EGI, USA). Central electrode E65 (Cz) was used as a reference. Four
364 electrodes were located above the eyes and on the cheeks to capture eye movements.
365 Electrode impedances were kept below 30 k Ω . The continuous EEG signal was
366 recorded throughout the session with a sampling rate of 500Hz.
367 Experiment sounds were played from a stereo speaker (Bose Corporation, USA)
368 positioned in front of the participant and at a 1m distance from the eyes; the sound level
369 was kept constant across participants and runs (70dB). The experiment ran on MATLAB
370 (R2018b; Natick, Massachusetts: The MathWorks Inc.); written instructions were
371 displayed only at the beginning of the experimental session, via Psychtoolbox version 3
372 (Brainard and Vision, 1997; PTB-3; <http://psychtoolbox.org/>).

373

374

375 **EEG Data Analysis**

376 **Preprocessing**

377 Data were preprocessed with a semi-automatic pipeline implemented in MATLAB (see
378 Stropahl et al., 2018; Bottari et al., 2020). Preprocessing was performed using EEGLAB
379 (Delorme and Makeig 2004; <https://scn.ucsd.edu/eeglab/index.php>). Data were loaded,
380 excluding electrode E65 (Cz), which was the reference channel of our EEG setup (thus
381 consisting only of zero values).

382 A high-pass filter (windowed sinc FIR filter, cut-off frequency 0.1 Hz, and filter order
383 10000) was applied to the continuous signal to remove slow drifts and DC offset.
384 A first segmentation in time was performed by epoching the signal according to the
385 event onset. To avoid boundary artifacts, the signal was cut 2s before its onset event
386 and until 2s after the end of the presentation (thus, from -2 to +114s) for each run. For
387 each participant, epochs were merged in a single file containing only the parts of the
388 signal referring to significant stimulation (thus excluding breaks between trials).
389 Independent Component Analysis (ICA; Bell and Sejnowski, 1995; Jung et al., 2000a,b)
390 was used to identify stereotypical artifacts. To improve ICA decomposition and reduce
391 computational time, data were low-pass filtered (windowed sinc FIR filter, cut-off
392 frequency 40Hz, filter order 50), downsampled to 250Hz, high-pass filtered (windowed
393 sinc FIR filter, cut-off frequency 1Hz, filter order 500), and segmented into consecutive
394 dummy epochs of 1s to spot non-stereotypical artifacts. Epochs with joint probability
395 larger than three standard deviations were rejected (Bottari et al., 2020). PCA rank
396 reduction was not applied before ICA to avoid compromising its quality and
397 effectiveness (Artoni, Delorme, and Makeig, 2018).
398 For each subject, ICA weights were computed using the EEGLAB runica algorithm and
399 then assigned to the corresponding original raw (unfiltered) dataset. Topographies for
400 each component were plotted for visual inspection. Artifacts associated with eye
401 movements and blinks were expected, and so a CORRMAP algorithm (Viola et al.,
402 2009) was used to remove components associated with such artifacts semi-
403 automatically. The automatic classification of components was performed using the
404 EEGLAB plugin ICLabel (Pion-Tonachini, Kreutz-Delgado, & Makeig, 2019).

405 Components representing eye movements and blinks were identified from their
406 topographical map within the components ICLabel marked as 'Eye' with a percentage
407 above 95%. Among these components, those with the highest rankings were selected
408 from a single dataset and used as templates (one for eye movements and one for
409 blinks). CORRMAP algorithm clusters ICA components with similar topography across
410 all datasets to highlight the similarity between the IC template and all the other ICs. A
411 correlation of the ICA inverse weights was computed, and similarity was allocated with a
412 threshold criterion of correlation coefficient being equal to or greater than 0.8 (default
413 value of CORRMAP; Viola et al., 2009). For all participants, on average, 1.92
414 components were removed (std= 0.88; range= 0-4).

415 Bad channels were interpolated after visually inspecting the scroll of the entire signal
416 and power spectral density for each electrode. On average, 3.75 (range= 1-8; std= 2.21)
417 channels were interpolated. The interpolation of noisy channels was performed via
418 spherical interpolation implemented in EEGLAB.

419 Finally, the reference channel (Cz) was reintroduced in the EEG data of each
420 participant, and the datasets were re-referenced to the average across all channels.
421

422 **Time domain analysis**

423 This analysis was performed to extract auditory evoked potentials and uncover phase-
424 locked magnitude changes associated with the two modes of sound representation
425 (Local Features or Summary Statistics).

426 Pre-processed data were low-pass filtered (windowed sinc FIR filter, cut-off frequency=
427 40Hz, filter order= 50). Additionally, detrend was applied by filtering the data above

428 0.5Hz (windowed sinc FIR filter, cut-off frequency= 0.5Hz, filter order= 2000).
429 Consecutive epochs (from -0.1 to 0.5s) were generated, including segments of either
430 the novel sounds or the repeated one (the second) of the triplets for each participant
431 and condition. Data were baseline corrected using the -0.1 to 0s pre-stimulus period.
432 Specifically, we averaged all the time points from -100 to 0ms before the onset of each
433 stimulus (either novel or repeated) and subtracted that value from post-stimulus activity
434 (Luck, 2014). Joint probability was used to prune non-stereotypical artifacts (i.e.,
435 sudden increment of muscular activation); the rejection threshold was four standard
436 deviations (Stropahl et al., 2018). For novel sounds, on average, 16.58 epochs per
437 participant were removed (std=5.42; range 5-30) out of the 144 concatenated epochs
438 that each Experiment * Duration comprised; for repeated sounds, on average, 16.15
439 epochs were removed (std= 5.11; range 5-29), again out of 144 trials per condition.
440 Data was converted from EEGLAB to FieldTrip (Oostenveld, Fries, Maris, and
441 Schoffelen, 2011; <http://fieldtriptoolbox.org>). Grand averages across participants were
442 computed for each experiment, duration, and stimulus type (repeated or novel). Data
443 across trials were averaged, generating Auditory Evoked Potentials (Figure 2-1 in
444 Extended Data).
445 For each triplet, we subtracted from the evoked response to the novel sound the one to
446 the preceding repeated one. Since all stimuli in the triplets (repeated and novel) were
447 never the same across runs and participants, the subtraction was performed to ensure
448 that neural responses were not driven by idiosyncratic differences in the stimuli that
449 were presented in that specific run, but by the statistical difference between novel and
450 repeated ones. Moreover, subtracting the response to the repeated sound from the one

451 to the novel sound allowed us to isolate within-triplet differences from those between
452 triplets. That is, since the first sound is repeated twice, the response to the second
453 repetition is not independent of the brain activity elicited by the first one and likely
454 incorporates a suppression mechanism to being exposed to the same stimulus twice. In
455 the same vein, the subtraction metrics represented the relative distance between being
456 exposed to the same sound as opposed to hearing a new one. Finally, the fact that in
457 the two experiments, novel and repeated sounds varied for selective properties (either
458 local features or summary statistics) allowed us to address how a deviation in fine
459 temporal details or global statistics altered the response to sound change.

460 A nonparametric permutation test was performed between experiments (Local Features
461 vs. Summary Statistics) for each duration (short, medium, and long), employing the
462 subtracted auditory responses between the novel and repeated sounds. The
463 permutation test was performed under the null hypothesis that probability distributions
464 across condition-specific averages were identical across experiments.

465 The cluster-based permutation approach is a nonparametric test that has the advantage
466 of solving the multiple comparison problem of multidimensional data in which you must
467 control several variables, such as time, space, frequencies, and experimental conditions
468 (Maris and Oostenveld, 2007).

469 Notably, statistical analyses between experiments were performed only within each
470 duration to avoid possible confounds associated with refractoriness effects due to
471 different interstimulus intervals (ISI) at long and short durations.

472 Thus, the contrasts of interest were: (1) Local Features short vs. Summary Statistics
473 short; (2) Local Features medium vs. Summary Statistics medium; (3) Local Features
474 long vs. Summary Statistics long.

475 A series of cluster-based permutation tests (Maris and Oostenveld, 2007; cluster alpha
476 threshold of 0.05 (two-tailed, accounting for positive and negative clusters); 10000
477 permutations; minimum neighboring channels = 2) was performed. Cluster-based
478 analyses were performed within a pool of central channels (according to EGI system,
479 channels: E3, E4, E6, E7, E9, E16, E21, E41, E51, E54, E65); we selected the
480 channels that better characterized the response to the second repeated sound, and
481 which corresponded to the 11 central sensors we used in the analysis (e.g., see the
482 topography in Extended Data, Figure 2-1). By pre-selecting this smaller number of
483 central channels (whose response likely originates from auditory sources), we avoided
484 including noisy channels in the model. Statistics were run for all samples from 0 to 0.5s.
485 We expected novel sounds to elicit larger responses than repeated sounds.

486

487 **Time-Frequency analysis**

488 Following the differences in magnitude changes observed between experiments for long
489 and short durations, we performed data decomposition in the time-frequency domain to
490 test whether sound changes at a high temporal resolution (local features in short
491 sounds) were encoded at faster timescales compared to those occurring at a low
492 temporal resolution (summary statistics in long sounds). We investigated frequencies
493 below 40Hz, which have been associated with auditory processing in studies including
494 both humans and animals (for review, see Gourevic et al., 2020). Specifically, several

495 studies have marked the relevance of lower (theta, alpha) and higher (beta) frequency
496 bands concerning auditory feature integration (e.g., VanRullen, 2016; Teng et al., 2018)
497 and detection of deviant sounds (e.g., Fujioka et al., 2012; Snyder and Large, 2005).
498 Preprocessed data were low-pass filtered to 100Hz (windowed sinc FIR filter, cut-off
499 frequency= 100Hz, filter order= 20) to attenuate high frequencies and high-pass filtered
500 at 0.5Hz (as with time-domain data). Data were epoched into segments from -0.5 to
501 1sec from stimulus onset: the second repeated or the novel. Joint probability was used
502 to remove bad segments with a threshold of 4 standard deviations. On average, 11.96
503 epochs were removed for repeated sounds (range= 4-25; std= 4.28) and 11.58 for novel
504 ones (range 4-26; std= 4.23). The resulting epoched datasets were converted to
505 Fieldtrip for time-frequency analysis. We used complex Morlet wavelets to extract the
506 power spectrum at each frequency of interest and time point. The frequencies spanned
507 from 4 to 40Hz in steps of 2Hz; the time window for decomposition comprised latencies
508 from -0.5 to 1s, around stimulus onset (either novel or repeated) in steps of 20ms.
509 Finally, the length of the wavelets (in cycles) increased linearly from 3 to 6.32 cycles
510 with increasing frequency (depending on the number of frequencies to estimate; N=19).
511 The signal was zero-padded at the beginning and end to ensure convolution with the
512 central part of the window. The resulting power spectrum for each participant was
513 averaged across trials.

514 Then, we performed a baseline correction to account for the power scaling ($1/f$). Unlike
515 ERP analysis, baseline selection is a more sensible choice in time-frequency.
516 Therefore, it was crucial to choose a baseline whose position did not affect the results
517 or over-boosted the effects. By using a stimulus-specific baseline as in the ERPs, for

518 the novel sounds, we would be using as baseline the activity from a condition that, at
519 least in some frequency ranges, is likely suppressed (the last 100ms of the response to
520 the second repeated sound), while for the second repeated sound, we would be using
521 as a baseline a segment in which activity is likely enhanced (as the first repeated sound
522 includes between-triplet changes). Because of the nonlinearity of the baseline (to
523 account for 1/f distribution), this will affect some frequencies more than others. When
524 subtracting the power to the second repeated sound from the power measured for the
525 novel sound, we would not be measuring the real dissimilarity between these
526 responses, because the baseline correction would be unfair and so the relative power
527 change. To account for this, we selected the same baseline for both the repeated and
528 novel sounds, corresponding to the activity from -100 to 0ms before the second
529 repeated sound. We decided to use a condition-averaged baseline (e.g., Cohen and
530 Donner, 2013; Cohen and Cavanagh, 2011) to account for differences in the oscillatory
531 tonic response as compared to the phasic one; since we are presenting a change
532 always at the same rate, the activity could be phase-locked in time in a similar way
533 across all the experiments, but the power at specific frequency bands could be higher in
534 one experiment as compared to the other. If we used a condition-specific baseline, this
535 effect would be masked because the activity would be corrected for the relative baseline
536 measured during that stimulation stream. Therefore, we took the activity from 100ms
537 prior to the onset of the second repeated sound for each experiment (Local Features or
538 Summary Statistics) and averaged their power separately for each duration. As a
539 baseline normalization method, we selected the relative change:
540 $(\text{pow}(t) - \text{bsl}) / \text{bsl}$

541 where pow is the total power at each sample (t) within the latencies of interest for
542 repeated and novel grand-averaged trials, and bsl is the averaged baseline (across
543 Experiment and time). The grand average of baseline-corrected power spectra of all
544 participants was computed.

545 We investigated the neural activity underlying the discrimination of novel and repeated
546 sounds across experiments for short and long durations. Thus, we first subtracted the
547 power at repeated trials from that at novel trials and then used cluster-based
548 permutation (Maris and Oostenveld, 2007) to investigate differences between neural
549 responses to sound changes across experiments (Local Features vs. Summary
550 Statistics) at each of the selected durations (short or long), at any latency (0-500ms)
551 and across all (65) channels (minimum neighboring channels = 1). Following the
552 inspection of power change between novel trials and repeated trials, oscillatory activity
553 above 30Hz was not considered. We used the period of the oscillatory activity as an
554 index of the temporal scale of the discriminative auditory processing, either slow,
555 medium or fast. Since we did not have any a priori hypothesis concerning the
556 contribution specific bands or ranges (e.g., from 9.5 to 16Hz) might have, we divided the
557 power change into equally spaced frequency bands (each including 8 frequencies of
558 interest, spaced in steps of 2Hz), creating a slow, medium, and fast oscillation range
559 between 4 and 30Hz. These frequencies of interest included canonical theta, alpha, and
560 beta oscillations (theta and alpha: 4-12Hz; low beta: 12-20Hz; high beta: 20-28Hz) but
561 were unbiased by their canonical subdivision (for which theta would be 4-7Hz, alpha 8-
562 12Hz, beta 13-25Hz and low gamma 25-40Hz). We instead hypothesized that the
563 temporal scale of oscillation (from slower to higher) would encode the type of change

564 that had occurred (local features vs. summary statistics). That is, depending on sound
565 duration, we expected to detect different power modulations in response to changes in
566 local features as compared to summary statistics at different timescales (frequency
567 bands). Cluster permutation was performed separately for each frequency range (10000
568 permutations). The directionality of the test was based on results in the Auditory Evoked
569 Responses (see Time-domain results) and on the specific frequency ranges:
570 specifically, for a short duration, we expected power changes in higher frequencies in
571 Local Features as compared to Summary Statistics. Conversely, at long duration, we
572 expected greater power changes in the lower-frequency range in response to sound
573 discrimination based on Summary Statistics compared with those based on Local
574 Features. For the short duration, we thus expected: Local Features > Summary
575 Statistics in the 4-12Hz range and Local Features < Summary Statistics in 12-20Hz and
576 20-28Hz. The opposite outcome was anticipated for the long duration: Summary
577 Statistics > Local Features in the alpha-theta range; Summary Statistics < Local
578 Features for beta bands (given the predefined directions of the effects, cluster alpha
579 threshold was 0.05, one-tailed).

580

581 **RESULTS**

582 **Behavioral Results**

583 For each condition, the percentage of correct beep detections was above 90% (Local
584 Features 40: mean= 0.99, std= 0.03; Local Features 209: mean= 0.99, std= 0.05; Local
585 Features 478: mean= 1, std= 0; Summary Statistics 40: mean=0.99, std= 0.05;
586 Summary Statistics 209: mean= 0.97, std=0.08; Summary Statistics 478: mean= 0.97,

587 std= 0.11; Figure 1-1A, in Extended Data). We ran a two-way ANOVA for repeated
588 measures with factors Experiment (2 levels, Local Features vs. Summary Statistics) and
589 Duration (3 levels, 40, 209, and 478) to address whether experiment type and stimulus
590 length had any impact on beep detection and participant attention to the task. No
591 significant main effects were observed (Experiment, $F(1,23)= 3.62$, $p = 0.07$, $\eta^2 = 0.14$;
592 Duration, $F(2,46)= 0.58$, $p = 0.56$, $\eta^2 = 0.3$) or their interaction (Experiment*Duration,
593 $F(2,46)= 0.45$, $p = 0.64$, $\eta^2 = 0.2$).

594 These behavioral results provide evidence that participants were attentive and
595 responsive during sound presentation throughout the experiment and that attention to
596 this orthogonal task was not influenced by the duration of the sound or experimental
597 condition.

598

599 **Time domain results**

600 By comparing Local Features vs. Summary Statistics separately for each sound
601 duration, cluster permutation revealed a significant positive cluster, selectively for the
602 short sound duration 40 ($p < 0.02$), lasting from 188 to 220ms after stimulus onset.
603 Following the prediction, results revealed a greater auditory potential of Local Features
604 compared to Summary Statistics for short duration. No significant positive cluster was
605 found for the medium (209) and long (478) sound durations (all $p > 0.39$). Conversely, a
606 significant negative cluster was found selectively for the long duration 478 ($p < 0.001$),
607 lasting from 220 to 308ms after stimulus onset. These results indicate a greater
608 response for Summary Statistics than Local features at long durations only. No
609 differences emerged for short and medium sound durations (all $ps > 0.33$).

610 Results clearly reveal double dissociations at the neural level based on stimulus length
611 and mode of representation (Figure 2B,C). Findings support behavioral outcomes for
612 which the processing of local features is favored for brief sound excerpts, while
613 summary statistics are built at a slower temporal rate as information is accumulated
614 (i.e., Berto et al., 2021; McDermott, Schemitsch, and Simoncelli, 2013). Going beyond
615 past behavioral effects, our results clearly show that local and summary representations
616 can emerge automatically from exposure to systematic sound changes. The neural
617 response to an acoustic change depends on the similarity between local features and
618 summary representations of sound excerpts. Summary statistics similarity can be
619 manipulated as a function of sound duration, eliciting a dissociation in the magnitude of
620 brain response that matches behavioral expectations.

621

622 **Time-Frequency Results**

623 Since summary statistics emerge over time, we expected statistical variations to be
624 encoded by slower oscillations than local feature changes. For such encoding, we
625 expected power modulations at faster oscillations in response to local feature changes
626 in short sounds and at slower oscillations in response to the emergence of a different
627 set of summary statistics in long acoustic excerpts. To test this, we separated the power
628 between 4 and 30Hz into three ranges, equally spaced: slow, 4-12Hz; medium, 16-
629 20Hz; and fast, 20-28Hz. Then, we used a nonparametric permutation approach to
630 address whether differences between Local Features and Summary Statistics emerged
631 according to sound duration (short or long) within the three frequency ranges.

632 Results followed the predicted pattern. For the short sound duration, the analysis
633 revealed a significant cluster between 100 and 220ms, in which sound change in Local
634 Features elicited a greater decrease of power in the fastest oscillation range (20-28Hz;
635 $p < 0.05$) compared to Summary Statistics (Figure 3A, left panel). This significant effect
636 was located over left frontocentral and right posterior sensors (see Grand-average
637 topography in Figure 3A, left). Conversely, for the long sound duration, we found a
638 greater increase of power in the slow oscillation range for Summary Statistics compared
639 to Local Features (4-12Hz; $p < 0.03$); the significant cluster consisted mainly of left
640 frontocentral channels and bilateral posterior channels and spanned from 260 to 500ms
641 (Figure 3A, right panel). No differences in power were found between Local Features
642 and Summary Statistics for any sound duration in the medium frequency range (12-
643 20Hz ranges, at any latency; all $ps > 0.24$). Overall, results revealed that when sound
644 duration is short, neural oscillations at higher frequency bands (canonically
645 corresponding to high-beta band) desynchronize more when the acoustic discrimination
646 is driven solely by local features; vice-versa when sound duration is long, i.e., higher
647 low-frequency oscillations (alpha and theta bands) are associated with stimulus
648 changes based on different summary statistics (Figure 3B).

649 Overall, these findings show that different temporal scales at the neural level underpin
650 the discrimination of variant elements in the auditory environment based on the amount
651 of information available and the type of sound change that has occurred.

652 Notably, beta desynchronization for Local Features (short duration) peaks 100-150ms
653 after stimulus onset, while the same effect in the time domain has a peak that builds up
654 around 200ms. The opposite was found for Summary Statistics (long duration), in which

655 theta-alpha synchronization starts about 40ms later than the effect observed in the time
656 domain and is more sustained over time (i.e., it lasts the entire time window). These
657 differences are indicative that the two measures capture at least partly different aspects
658 of sound discrimination.

659

660 **DISCUSSION**

661 The auditory system extracts information at high (local) and low (summary) temporal
662 resolutions. Here, we assessed whether discriminative responses to local or summary
663 representations could be measured at the neural level and whether they are encoded at
664 different temporal scales (Panzeri et al., 2010). We employed a computational model
665 (McDermott and Simoncelli, 2011) to synthetically create stimuli with the same summary
666 statistics but different local features. We used these synthetic stimuli to present streams
667 of triplets containing repeated and novel sounds that could vary in their local features or
668 summary statistics.

669 Results in the time domain showed that when the sound duration was short, the
670 magnitude of auditory potentials increased selectively for changes in local features. In
671 contrast, when the sound duration was long, changes in auditory statistics elicited a
672 higher response compared with changes in local features (Figure 2B, C). Thus,
673 according to sound duration, we observed an opposite trend in the magnitude change of
674 the evoked response. Note that for each sound texture, we manipulated the duration of
675 the excerpts, and not their properties (we synthesized only 4 synthetic exemplars per
676 sound texture, that we cut into smaller excerpts either 40, 209, or 478ms which were
677 then randomly drawn in the experiments; see Material and Methods above). The

678 dissociation observed between experiments according to sound duration is indicative
679 that the processing of local features or summary statistics is strictly dependent on the
680 amount of information presented. This trend perfectly matched expectations based on
681 previous psychophysics evaluations (i.e., Berto et al., 2021; McDermott et al., 2013)
682 despite the protocol was slightly different from the behavioral implementation. In the
683 psychophysical version, the two experiments (Local and Summary) were substantially
684 different from each other. One experiment, called Exemplar Discrimination, was the
685 equivalent of the Local Features experiment in our protocol and contained two different
686 sounds (since one was repeated twice). However, the other experiment, named Texture
687 Discrimination, contained three different sound excerpts (two derived from the same
688 white noise but with different imposed summary statistics; one derived from a different
689 white noise with the same statistics). Different task demands justified this disparity.
690 Specifically, in the behavioral version, participants were given very clear instructions on
691 which sound properties to pay attention to during each experiment (sound details or
692 sound source, respectively) and even which sound to use for comparison (the middle
693 one; McDermott et al., 2013). In this protocol, the sequences had the same structure in
694 both experiments (two repeated sounds followed by a novel one), while the only
695 difference was the generative statistics imposed on the novel sound compared to the
696 repeated one (same in Local and different in Summary) or the white noise sample used
697 to initialize the synthesis (different in Local and same in Summary). This allowed us to
698 test for the automaticity of the processes and to measure distinct neural responses
699 when the system is exposed to a similar or different set of statistics combined with the
700 same or different local structure. Moreover, it permitted a fair comparison between

701 experiments. Nonetheless, results went in the same direction in both the EEG and the
702 behavioral evaluations, suggesting similar mechanisms are in place despite the lack of
703 an explicit request to pay attention to specific sound properties.
704 Finally, analysis in the time-frequency domain revealed that neural activity at different
705 temporal scales characterized discriminative responses to local features or summary
706 statistics. Faster oscillations (in the beta range) were associated with discrimination
707 based on local features, and slower oscillations (in the theta-alpha range) with changes
708 based on summary statistics.

709

710 **Automaticity of Local Features and Summary Statistics Processing**

711 Auditory responses to novel local features or summary statistics were associated with
712 differences in magnitude that could be automatically detected. This finding confirms that
713 the auditory system can attune its response to specific sound changes and expands
714 seminal studies measuring the mismatch negativity (MMN) response (Näätänen et al.,
715 1978; Tiitinen et al., 1994). MMN is the neural marker of a process by which the system
716 “scans” for regularities in entering sounds and uses them as references to detect
717 variations in the auditory scene (for reviews, see Näätänen et al., 2001, 2010). In our
718 study, expectations that a change would occur in the third element of the triplet had a
719 probability of 1 in each experiment (Local Features and Summary Statistics; Figure 2A).
720 Thus, spurious expectancy or attentional effects cannot explain results. Coherently, the
721 MMN response to a deviant sound is not affected by prior expectations that the novel
722 element will occur (Rinne et al., 2001); rather, the auditory system automatically orients
723 attention toward it. Here we highlighted another ability of the system. Beyond automatic

724 orientation toward a relevant deviant sound, our results show that it is possible to
725 categorize the acoustic change according to the representation (local or summary) and
726 temporal resolution (high or low) at which it has occurred. Importantly, discriminative
727 neural responses could be detected even if the task per se did not involve any
728 discrimination or in-depth processing of either local features or summary statistics. In
729 other words, the sound changes were processed even when irrelevant to the behavioral
730 task participants were attending (rare beep detection), strongly suggesting that the
731 entrainment to local or global acoustic change emerges automatically from exposure to
732 regular changes in the environment and is strictly dependent on the amount of
733 information presented.

734 Furthermore, the double dissociation we observed based on sound duration (with Local
735 Features eliciting greater magnitude change than Summary Statistics for short sounds
736 and vice-versa for long sounds) rules out the possibility of results being explained by a
737 mere saliency effect (i.e., the fact that, in Summary Statistics, a different sound object
738 was presented). Importantly, the main advantage of using synthetic sounds instead of
739 natural recordings was to be able to control the summary statistics embedded in the
740 sounds. That is, all sounds were random white noise samples to which we imposed the
741 same (or a different set) of summary statistics. If the brain were not automatically
742 encoding the summary statistics, we would not have been able to distinguish between
743 Local Features and Summary Statistics experiments, especially at long duration, since
744 all repeated and novel sounds differed for their local structure. Nor would it have been
745 possible to detect a dissociation in the neural response according to sound duration.
746 This observation is further supported by the fact that results emerged despite sound

747 objects between the triplets being continuously changing (the only fixed parameter was
748 the expected similarity in local features or summary statistics between the novel and
749 repeated sounds).

750 These findings can be generalized to a variety of sound textures (Figure 2A; see also
751 Extended Data, Table 1-2) and the exact moment in which the summary percepts
752 emerge likely depends on specific comparisons across sound objects (repeated and
753 novel). In line with this, using many different sounds to create sound streams led to
754 grand averaged signals associated with discrimination based on summary statistics with
755 a rather spread-out shape (see Figure 2C, right).

756 Finally, it is important to notice that imposing different statistics on the same white noise
757 leads to sounds with different long-term average spectra. Therefore, it is possible that
758 magnitude differences in response to the Summary Statistics experiment, compared to
759 Local Features, were driven by low-level spectrotemporal modulations rather than
760 changes in higher-order statistics. However, if that was the case, we might have
761 expected an effect already at medium duration (209ms), which was instead not present.
762 Further experiments may be required to fully rule out this possible confound.

763

764 **Local features changes are encoded by fast oscillations**

765 By comparing the difference in total power between novel and repeated sounds in the
766 two experiments, we found that, for short sounds, the power between 20 and 28Hz
767 decreased when a change in local features was detected, as compared to when
768 summary statistics were changed. This desynchronization occurred between 80 and
769 200ms after stimulus onset (Figure 3A, B, left). Desynchronization of oscillatory activity

770 is the decrease in power measured at specific frequency bands (generally alpha and
771 beta ranges), which generally emerges following the onset of an event (Pfurtscheller
772 and Lopes da Silva, 1999). It results from increased cellular excitability in thalamocortical
773 circuits and generally reflects cortical withdrawal from the resting state to engage in a
774 cognitive process (Pfurtscheller and Lopes da Silva, 1999).

775 The 20-28Hz band includes frequencies that are canonically attributed to high-beta
776 oscillations. Changes in power synchronization in the beta range have been correlated
777 with performance in tasks involving the detection of temporal or intensity deviations
778 (Arnal et al., 2015; Herrmann et al., 2016). Overall, these findings suggest that, among
779 other operations, brain activity in the high beta range could be engaged in the
780 processing of low-level properties of a stimulus. Beta-band activity has also been
781 investigated in the context of rhythmic perception. A disruption in beta power can be
782 observed in non-rhythmic sequences or when an attended tone is omitted from a
783 regular series (e.g., Fujioka et al., 2012). Interestingly, beta synchronization not only
784 captures irregularities in a pattern but also reflects the type of change that has occurred.
785 For instance, it has been shown that beta desynchronization was higher prior to the
786 occurrence of a deviant sound whose pitch varied in a predictable way, as compared to
787 an unpredictable variation. Accordingly, beta desynchronization has been proposed as
788 a marker of predictive coding (Engel and Fries, 2010; Chang, Bosnyak, and Trainor,
789 2018).

790 In our model, stimuli could be derived from the same white noise sample or a different
791 one (Figure 1C). In Local Features, the novel sound is derived from another white noise
792 sample, as compared to the repeated sound on which we imposed the same summary

793 statistics. Thus, with this synthesis approach, in terms of fine acoustic features, when
794 sounds were short, novel sounds had a more different temporal structure (Figure 1F)
795 and were statistically more dissimilar (Figure 1D, 2B) than their paired repeated one in
796 the Local Features experiment as compared to Summary Statistics. Overall, these
797 results suggest that, in the absence of enough information to build summary
798 representations, faster oscillations are in charge of small, acoustic change detection to
799 be used to discriminate sound excerpts.

800

801 **Slower oscillations are engaged in Summary Statistics processing**

802 By comparing Local Features with Summary Statistics at long durations, we observed
803 that the emergence of different auditory statistics in the novel sound, as compared to
804 the previous, repeated one, elicited higher power at slower frequencies, compatible with
805 canonical alpha-theta oscillations. This power synchronization emerged at relatively late
806 latencies from stimulus onset (between 240 and 500ms; Figure 3A, B, right). and was
807 not present when solely local features were driving sound change (as in the Local
808 Features experiment). Provided that summary statistics can primarily be measured at
809 increasing sound duration, we expected differences between long-duration stimuli being
810 carried by relatively slower brain activity. However, statistical comparisons were
811 performed within the sound duration; thus, if this effect was simply driven by the sounds
812 being longer (478ms) rather than the processing of auditory statistics, we should not
813 have observed a difference in alpha-theta synchronization between experiments.
814 Similarly, if the effect were driven by simply presenting a “more different” sound in the
815 Summary Statistics experiment, as compared to Local Features one, then we would

816 have seen an effect also for 209ms, which was not the case; similarly, we would not
817 have been able to dissociate the effects based on sound duration. Finally, it is worth
818 noting that the stimulation rate was kept constant across all tested durations (40, 209,
819 and 478), meaning that we always presented one sound every half a second. This
820 means that, disregarding the amount of information we presented, the change always
821 occurred in a window of 1.5 seconds (with novel sound always occurring at a frequency
822 of 0.667Hz). Therefore, the effect strictly depends on the amount of information we
823 presented within this temporal window, rather than the time interval between sound
824 excerpts.

825 A previous study investigated the temporal window of integration of sound textures,
826 showing that it can extend for several seconds (McWalter and McDermott, 2018, 2019).
827 In this study, we could not use stimuli longer than 500ms to maintain the 2Hz rhythmic
828 stimulation pattern in all experiments. Thus, we could not address the integration effects
829 of single sounds at longer durations. Interestingly, the integration window measured for
830 sound textures is relatively long compared to the receptive fields of auditory neurons,
831 whose response has been shown to be sustained for about a few hundred milliseconds
832 (e.g., Miller et al., 2002). Previous evidence suggested the existence of an active
833 chunking mechanism condensing entering acoustic information within a much longer
834 temporal window, approximately 150-300ms (VanRullen, 2016; Riecke, Sacks, &
835 Schroeder., 2015; Teng et al., 2018). Such integration length would be related to
836 ongoing oscillatory cycles, specifically corresponding to the theta range (4-7Hz; Ghitza
837 & Greenberg, 2009; Ghitza, 2012). Compatibly, a recent study showed that acoustic

838 changes occurring around 200ms could explain the modulations of phase
839 synchronization in theta (Teng et al., 2018).
840 Although there is no evidence that 200ms windows are relevant for texture perception
841 (see McWalter and McDermott, 2018, 2019), our data show that brain activity already
842 synchronizes 200ms after stimulus onset to the emergence of a novel set of auditory
843 statistics. The integration window of sound texture defined by previous studies refers to
844 the maximum duration within which the averaging of local information into summary
845 statistics can occur (McWalter and McDermott, 2018). It is still unclear how this relates
846 to the emergence of relevant percepts in the brain (i.e., sound object identity) in
847 response to average statistics. The higher power synchronization in the theta-alpha
848 range observed in response to sensory statistics might be interpreted as one of the
849 possible neural mechanisms underlying the development of such abstract
850 representations, which may lead to the perceptual understanding that a new sound
851 object has occurred. This would explain why it happens when a different set of statistics
852 is detected and not when only local features change while sound identity remains
853 unchanged.

854

855 **CONCLUSION**

856 Combining a computational synthesis approach with electrophysiology revealed distinct
857 cortical representations associated with local and summary representations. We
858 showed that different neural codes at faster and slower temporal scales are entrained to
859 automatically detect changes in entering sounds based on summary statistics similarity
860 emerging as a function of sound duration. These results promote using computational

861 methods to appoint neural markers for basic auditory computation in fundamental and
862 applied research. Furthermore, the automaticity of the protocol and the fast
863 implementation allow the testing of different populations (including newborns, infants,
864 children, and clinical patients) that do not have the resources to attend to complex
865 tasks.

866

867 **DATA AVAILABILITY**

868 Raw EEG data, analysis scripts, participants' information, and sound excerpts employed
869 in the experiment are available in an online repository at this link:

870 <https://data.mendeley.com/datasets/gx7cb7fzv4/1>

871

872 **REFERENCES**

873 McAdams, S. (1993). Recognition of sound sources and events. *Thinking in sound: The*
874 *cognitive psychology of human audition*, 146-198.

875

876 Griffiths, T. D. (2001). The neural processing of complex sounds. *Annals of the New York*
877 *Academy of Sciences*, 930(1), 133-14

878

879 Plomp, R. (1964). Rate of decay of auditory sensation. *The Journal of the Acoustical Society of*
880 *America*, 36(2), 277-282.

881

882 McDermott, J. H., Schemitsch, M., & Simoncelli, E. P. (2013). Summary statistics in auditory
883 perception. *Nature Neuroscience*, 16(4), 493-498.

884

885 McDermott, J. H., & Simoncelli, E. P. (2011). Sound texture perception via statistics of the
886 auditory periphery: evidence from sound synthesis. *Neuron*, 71(5), 926-940.
887

888 Dau, T., Kollmeier, B., & Kohlrausch, A. (1997). Modeling auditory processing of amplitude
889 modulation. I. Detection and masking with narrow-band carriers. *The Journal of the Acoustical*
890 *Society of America*, 102(5), 2892-2905.
891

892 Yabe, H., Tervaniemi, M., Sinkkonen, J., Huotilainen, M., Ilmoniemi, R. J., & Näätänen, R.
893 (1998). Temporal window of integration of auditory information in the human brain.
894 *Psychophysiology*, 35(5), 615-619.
895

896 Saint-Arnaud, N., Popat, K. (2021). Analysis and synthesis of sound textures
897 D.F. Rosenthal, H.G. Okuno (Eds.), *Computational Auditory Scene Analysis*, CRC Press (2021),
898 pp. 293-308
899

900 Barlow, H. B. (1961). Possible principles underlying the transformation of sensory messages
901 *Sensory Communication* ed WA Rosenblith. Cambridge, MA: MIT Press), 7, 1-34.
902

903 Berto, M., Ricciardi, E., Pietrini, P., & Bottari, D. (2021). Interactions between auditory statistics
904 processing and visual experience emerge only in late development. *Iscience*, 24(11), 103383.
905

906 Zhai, X., Khatami, F., Sadeghi, M., He, F., Read, H. L., Stevenson, I. H., & Escabí, M. A. (2020).
907 Distinct neural ensemble response statistics are associated with recognition and discrimination
908 of natural sound textures. *Proceedings of the National Academy of Sciences*, 117(49), 31482-
909 31493.
910

- 911 Panzeri, S., Brunel, N., Logothetis, N. K., & Kayser, C. (2010). Sensory neural codes using
912 multiplexed temporal scales. *Trends in neurosciences*, 33(3), 111-120.
913
- 914 Wang, C., & Zhang, Q. (2021). Word frequency effect in written production: Evidence from
915 ERPs and neural oscillations. *Psychophysiology*, 58(5), e13775.
916
- 917 World Medical Association. (2013). World Medical Association Declaration of Helsinki: ethical
918 principles for medical research involving human subjects. *Jama*, 310(20), 2191-2194.
919
- 920 Bacon, S. P., & Grantham, D. W. (1989). Modulation masking: Effects of modulation frequency,
921 depth, and phase. *The Journal of the Acoustical Society of America*, 85(6), 2575-2580.
922
- 923 Sussman, E. S., & Gumenyuk, V. (2005). Organization of sequential sounds in auditory
924 memory. *Neuroreport*, 16(13), 1519-1523.
925
- 926 Brainard, D. H., & Vision, S. (1997). The psychophysics toolbox. *Spatial vision*, 10(4), 433-436.
927
- 928 Stropahl, M., Bauer, A. K. R., Debener, S., & Bleichner, M. G. (2018). Source-Modeling auditory
929 processes of EEG data using EEGLAB and brainstorm. *Frontiers in neuroscience*, 12, 309.
930
- 931 Bottari, D., Bednaya, E., Dormal, G., Villwock, A., Dzhelyova, M., Grin, K., ... & Röder, B.
932 (2020). EEG frequency-tagging demonstrates increased left hemispheric involvement and
933 crossmodal plasticity for face processing in congenitally deaf signers. *NeuroImage*, 223,
934 117315.
935

- 936 Delorme, A., & Makeig, S. (2004). EEGLAB: an open source toolbox for analysis of single-trial
937 EEG dynamics including independent component analysis. *Journal of neuroscience methods*,
938 134(1), 9-21.
- 939
- 940 Bell, A. J., & Sejnowski, T. J. (1995). An information-maximization approach to blind separation
941 and blind deconvolution. *Neural computation*, 7(6), 1129-1159.
- 942
- 943 Jung, T. P., Makeig, S., Humphries, C., Lee, T. W., Mckeown, M. J., Iragui, V., & Sejnowski, T.
944 J. (2000a). Removing electroencephalographic artifacts by blind source separation.
945 *Psychophysiology*, 37(2), 163-178.
- 946
- 947 Jung, T. P., Makeig, S., Westerfield, M., Townsend, J., Courchesne, E., & Sejnowski, T. J.
948 (2000b). Removal of eye activity artifacts from visual event-related potentials in normal and
949 clinical subjects. *Clinical Neurophysiology*, 111(10), 1745-1758.
- 950
- 951 Artoni, F., Delorme, A., & Makeig, S. (2018). Applying dimension reduction to EEG data by
952 Principal Component Analysis reduces the quality of its subsequent Independent Component
953 decomposition. *NeuroImage*, 175, 176-187.
- 954
- 955 Viola, F. C., Thorne, J., Edmonds, B., Schneider, T., Eichele, T., & Debener, S. (2009). Semi-
956 automatic identification of independent components representing EEG artifact. *Clinical*
957 *Neurophysiology*, 120(5), 868-877.
- 958
- 959 Pion-Tonachini, L., Kreutz-Delgado, K., & Makeig, S. (2019). ICLabel: An automated
960 electroencephalographic independent component classifier, dataset, and website. *NeuroImage*,
961 198, 181-197.

962

963 Luck, S. J. (2014). An introduction to the event-related potential technique. MIT press.

964

965 Oostenveld, R., Fries, P., Maris, E., & Schoffelen, J. M. (2011). FieldTrip: open source software
966 for advanced analysis of MEG, EEG, and invasive electrophysiological data. *Computational
967 intelligence and neuroscience*, 2011.

968

969 Maris, E., & Oostenveld, R. (2007). Nonparametric statistical testing of EEG-and MEG-data.
970 *Journal of neuroscience methods*, 164(1), 177-190.

971

972 Gourévitch, B., Martin, C., Postal, O., & Eggermont, J. J. (2020). Oscillations in the auditory
973 system and their possible role. *Neuroscience & Biobehavioral Reviews*, 113, 507-528.

974

975 VanRullen, R. (2016). Perceptual cycles. *Trends in cognitive sciences*, 20(10), 723-735.

976

977 Teng, X., Tian, X., Doelling, K., & Poeppel, D. (2018). Theta band oscillations reflect more than
978 entrainment: behavioral and neural evidence demonstrates an active chunking process.

979 *European Journal of Neuroscience*, 48(8), 2770-2782.

980

981 Fujioka, T., Trainor, L. J., Large, E. W., & Ross, B. (2012). Internalized timing of isochronous
982 sounds is represented in neuromagnetic beta oscillations. *Journal of Neuroscience*, 32(5), 1791-
983 1802.

984

985 Snyder, J. S., & Large, E. W. (2005). Gamma-band activity reflects the metric structure of
986 rhythmic tone sequences. *Cognitive brain research*, 24(1), 117-126.

987

- 988 Cohen, M. X., & Donner, T. H. (2013). Midfrontal conflict-related theta-band power reflects
989 neural oscillations that predict behavior. *Journal of neurophysiology*, 110(12), 2752-2763.
990
- 991 Cohen, M. X., & Cavanagh, J. F. (2011). Single-trial regression elucidates the role of prefrontal
992 theta oscillations in response conflict. *Frontiers in psychology*, 2, 30.
993
- 994 Näätänen, R., Gaillard, A. W., & Mäntysalo, S. (1978). Early selective-attention effect on evoked
995 potential reinterpreted. *Acta psychologica*, 42(4), 313-329.;
996
- 997 Tiitinen, H., May, P., Reinikainen, K., & Näätänen, R. (1994). Attentive novelty detection in
998 humans is governed by pre-attentive sensory memory. *Nature*, 372(6501), 90-92.
999
- 1000 Näätänen, R., Tervaniemi, M., Sussman, E., Paavilainen, P., & Winkler, I. (2001). 'Primitive
1001 intelligence' in the auditory cortex. *Trends in neurosciences*, 24(5), 283-288.
1002
- 1003 Näätänen, R., Astikainen, P., Ruusuvirta, T., & Huotilainen, M. (2010). Automatic auditory
1004 intelligence: An expression of the sensory–cognitive core of cognitive processes. *Brain research*
1005 *reviews*, 64(1), 123-136.
1006
- 1007 Rinne, T., Antila, S., & Winkler, I. (2001). Mismatch negativity is unaffected by top-down
1008 predictive information. *NeuroReport*, 12(10), 2209-2213.
1009
- 1010 Pfurtscheller, G., & Da Silva, F. L. (1999). Event-related EEG/MEG synchronization and
1011 desynchronization: basic principles. *Clinical neurophysiology*, 110(11), 1842-1857.
1012

- 1013 Arnal, L. H., Doelling, K. B., & Poeppel, D. (2015). Delta–beta coupled oscillations underlie
1014 temporal prediction accuracy. *Cerebral Cortex*, 25(9), 3077-3085.
1015
- 1016 Herrmann, C. S., Strüber, D., Helfrich, R. F., & Engel, A. K. (2016). EEG oscillations: from
1017 correlation to causality. *International Journal of Psychophysiology*, 103, 12-21.
1018
- 1019 Engel, A. K., & Fries, P. (2010). Beta-band oscillations—signalling the status quo?. *Current*
1020 *opinion in neurobiology*, 20(2), 156-165.
1021
- 1022 Chang, A., Bosnyak, D. J., & Trainor, L. J. (2018). Beta oscillatory power modulation reflects the
1023 predictability of pitch change. *Cortex*, 106, 248-260.
1024
- 1025 McWalter, R., & McDermott, J. H. (2018). Adaptive and selective time averaging of auditory
1026 scenes. *Current Biology*, 28(9), 1405-1418.
1027
- 1028 McWalter, R., & McDermott, J. H. (2019). Illusory sound texture reveals multi-second statistical
1029 completion in auditory scene analysis. *Nature communications*, 10(1), 5096.
1030
- 1031 Miller, L. M., Escabi, M. A., Read, H. L., & Schreiner, C. E. (2002). Spectrotemporal receptive
1032 fields in the lemniscal auditory thalamus and cortex. *Journal of neurophysiology*, 87(1), 516-527.
1033
- 1034 Riecke, L., Sack, A. T., & Schroeder, C. E. (2015). Endogenous delta/theta sound-brain phase
1035 entrainment accelerates the buildup of auditory streaming. *Current Biology*, 25(24), 3196-3201.
1036
- 1037 Viemeister, N. F. (1979). Temporal modulation transfer functions based upon modulation
1038 thresholds. *The Journal of the Acoustical Society of America*, 66(5), 1364-1380.

1039

1040 Ghitza, O., & Greenberg, S. (2009). On the possible role of brain rhythms in speech perception:
1041 intelligibility of time-compressed speech with periodic and aperiodic insertions of silence.
1042 *Phonetica*, 66(1-2), 113-126.

1043

1044 Ghitza, O. (2012). On the role of theta-driven syllabic parsing in decoding speech: intelligibility
1045 of speech with a manipulated modulation spectrum. *Frontiers in psychology*, 3, 238

1046

1047

1048 **Legends**

1049 **Figure 1. Experimental stimuli.** (A) Computational Texture Model to extract auditory
1050 statistics. An original recording of a natural sound texture is passed through the auditory
1051 texture model (the list of presented sound textures is available as Extended Data,
1052 Figure 1-2). The model provides a mathematical formulation of the auditory system's
1053 computations (auditory statistics) to represent the sound object. The signal is filtered
1054 with 32 audio filters to extract analytic and envelope modulations for each cochlear sub-
1055 band. Envelopes are downsampled and multiplied by a compression factor. From the
1056 compressed envelopes, a first set of statistics is computed: marginal moments
1057 (including envelope mean, variance, and skewness), autocorrelation between temporal
1058 intervals, and cross-band correlations. Compressed envelopes are then filtered with 20
1059 modulation filters. The remaining statistics are extracted from the filtered envelopes:
1060 modulation power and cross-band correlations between envelopes filtered with the
1061 same modulation filter (C1) and between the same envelope filtered through different
1062 filters (C2).

1063 (B) Schematic of Sound Synthesis. The white-noise sample is filtered through the
1064 auditory model (McDermott and Simoncelli, 2011) to extract its cochlear envelopes,
1065 which are then subtracted from those obtained from the original sound texture. The
1066 average statistics from the original sound textures are then imposed on the subtracted
1067 white noise envelopes. The outcome is multiplied by the fine structure of the white noise
1068 sample to preserve its local acoustic distribution (e.g., temporal structure). The result is
1069 recombined in the synthetic signal, reiterating the procedure until a desired SNR of 20-
1070 dB is reached.

1071 (C) Impact of white noise sample and imposed statistics on synthetic sounds. Two
1072 different sets of statistics are extracted from two sound textures: “frogs” and “horse
1073 trotting”. Each set of values is imposed on two different random white noise samples.
1074 When the same statistics are imposed on different white noise samples, the outcomes
1075 are two synthetic exemplars of the same sound texture. These exemplars will have the
1076 same summary statistical representation but will diverge in their local features as the
1077 original input sound will influence them. When different statistics are imposed on the
1078 same white noise sample, the results are two synthetic exemplars that will diverge in
1079 their overall summary statistics and be perceptually associated with different sound
1080 objects. The cochleograms of the 0.5 s synthetic exemplars are displayed.

1081 (D) Similarity of statistics between excerpt pairs. Couples of sound excerpts presented
1082 in the study (repeated and novel; see Figure 2A for the experimental protocol) could be
1083 derived from different white noise samples to which we imposed the same statistics (in
1084 coral) or from the same white noise sample with different statistics (in blue). The
1085 summary statistics similarity between these couples of synthetic excerpts was

1086 computed by averaging the SNRs between statistics of repeated and novel sounds,
1087 measured separately for each statistical class. Boxplots show the averaged SNRs at
1088 three sound durations of interest (short, 40ms; medium, 209ms; long, 478ms). When
1089 sounds were short (40ms), statistical values were more similar for sounds derived from
1090 the same white noise samples (in blue) compared to different ones (in coral), even
1091 when including different original statistics. As duration increased (209, 478ms), statistics
1092 progressively converged to their original values and were more dissimilar for sounds
1093 with different generative statistics (blue) than for sounds including the same statistics
1094 (coral), irrespectively of original white noise sample. *** $p < 0.001$

1095 (E) Comparing auditory statistics of 478ms synthetic sounds. Envelope marginal
1096 moments (mean, skewness, and variance) of all sound textures are displayed, while
1097 highlighted are those from three sound excerpts selected randomly; two have the same
1098 imposed auditory statistics (in red and yellow), and one has different statistics (in blue).
1099 In the bottom row, the remaining statistics are displayed (envelope correlation,
1100 modulation power, C1, and C2). The similarity between statistical values is higher when
1101 the sounds come from the same original texture.

1102 (F) Similarity between envelope pairs of short sounds. In the top panel, boxplots
1103 represent the correlation coefficients (r) measured between broadband envelopes for
1104 each pair of 40ms sound excerpts (repeated and novel; $n = 6912$) divided according to
1105 experiment (Local Features or Summary Statistics). Amplitude modulations of brief
1106 excerpts are significantly more similar when sound pairs originate from the same white
1107 noise sample (Summary Statistics experiment) than when they do not (as in the Local
1108 Features experiment), disregarding their imposed generative statistics. *** $p < 0.001$.

1109 In the bottom panel, show examples of the 40ms broadband envelopes used for
1110 computing the correlation coefficients (r) above.

1111

1112 **Figure 2. Experimental procedure and results of time domain analysis. (A)**

1113 Experimental protocol for EEG. Triplets of sounds were presented at a fast rate (one
1114 sound every 500ms). Two sounds were identical (Repeated), while the third was
1115 different (Novel) and could vary in its local features (left) or summary statistics (right)
1116 depending on the experiment (Local Features or Summary Statistics). Three sound
1117 durations, equally spaced logarithmically (short, medium, and Long: 40, 209, and
1118 478ms), were employed (in different sound streams) to tap into each auditory mode
1119 separately (local features vs. summary statistics processing). The list of presented
1120 sound textures is available as Extended Data, Figure 1-2. To ensure participants were
1121 attentive during the presentation, they performed an orthogonal task, consisting of
1122 pressing a button when an infrequent target (beep) appears. Performance accuracy was
1123 high in all experiments and durations and is displayed in Figure 1-1 in Extended Data.

1124 (B) Grand average topographies of the differential response associated with the sound
1125 change (novel sound minus repeated sound) at significant latencies for each experiment
1126 and duration. For each latency, electrodes associated with significant clusters are
1127 displayed above as red stars on the scalp. * $p < 0.025$.

1128 On the right side of the topographical maps, the boxplots represent objective differences
1129 between the novel and repeated sounds of all auditory statistics (averaged). The
1130 difference was computed between the statistics of sounds presented for each run,
1131 experiment, and duration and averaged across all participants. Within each duration,

1132 medians differed at the 5% significance level between experiments. Local Features >
1133 Summary Statistics at short (40) duration and Summary Statistics > Local Features for
1134 medium (209) and long (478) durations. The evoked response in the EEG agrees with
1135 the objective statistical difference measured from the sound excerpts.

1136 (C) Grand average electrical activity (negative values are plotted up) of the differential
1137 response (novel minus repeated) at significant electrodes (in red) for both short and
1138 long durations. Shaded regions show interpolated repeated error of the mean (SE) at
1139 each time point. Positive values indicate that the novel elicited a greater response than
1140 repeated. Results of cluster permutation are displayed as black bars extending through
1141 significant latencies. – $p < 0.025$.

1142 For visualizing the ERPs before subtraction (novel -repeated), see Extended Data,
1143 Figure 2-1.

1144

1145 **Figure 3. Results of time-frequency analysis.** (A) Grand average difference (novel
1146 minus repeated) of total power for short and long sound durations in both experiments
1147 (Local Features and Summary Statistics) at significant channels. Rectangular regions
1148 comprise the latencies and frequency range in which power changes were significant
1149 between experiments after cluster-based permutation. Significant channels are marked
1150 as red stars over the sketch of a scalp ($* p < 0.05$). In the left panel, results for the short
1151 duration are displayed and show higher-power desynchronization in the 20-28Hz
1152 frequency range (high beta) for Local Features as compared to Summary Statistics. In
1153 the right panel, results for the long duration show higher 4-12Hz (alpha-theta) power
1154 synchronization for Summary Statistics as compared to Local Features. Grand-average

1155 topographical maps at significant latencies and frequency ranges are displayed next to
1156 the corresponding power-spectrum plots.

1157 (B) Average power difference between novel and repeated sounds for each range of
1158 frequency bands (Slow, Medium, and Fast), averaged across all significant channels,
1159 plotted at all latencies (from 0 to 0.5s). Significant channels are marked as red stars
1160 over the sketch of a scalp. Shaded regions show interpolated standard error of the
1161 mean (SE) at each time point. * $p < 0.05$.

1162

1163 **Figure 1-1. Behavioral results. Related to Figure 2.** (A) The group-level average
1164 proportion of correct detections of beeps when presented. Bar plots represent average
1165 values of hits across all participants. Error bars represent the standard error of the
1166 mean (SE). No significant difference existed across conditions (all $p > 0.05$).

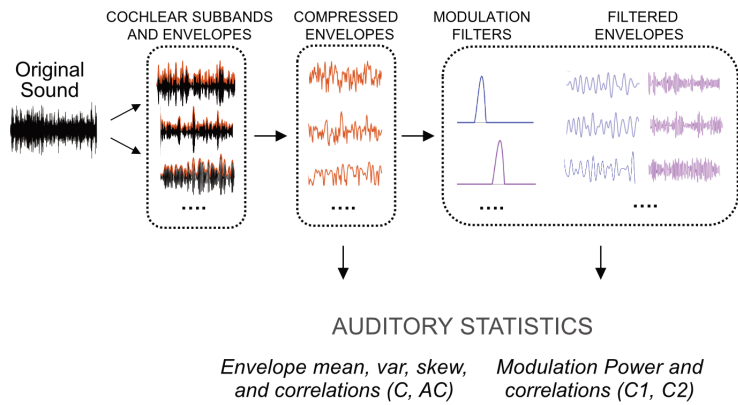
1167

1168 **Figure 1-2. List of Sound Textures. Related to Figure 1 and 2.** In Local Features
1169 Discrimination, for each sound texture in column 1, two synthetic exemplars of the
1170 sound texture were selected. One was presented twice (repeated) and the other was
1171 presented as the third element of the triplet (novel). In Summary Statistics
1172 Discrimination, sound textures were paired according to perceived similarity
1173 (McDermott, Schemitsch, and Simoncelli, 2013). For each sound texture in column 1,
1174 one synthetic exemplar was selected and presented twice. Then, an exemplar of the
1175 texture from the corresponding row in column 2 was selected and used as the third
1176 element of the triplet (novel).

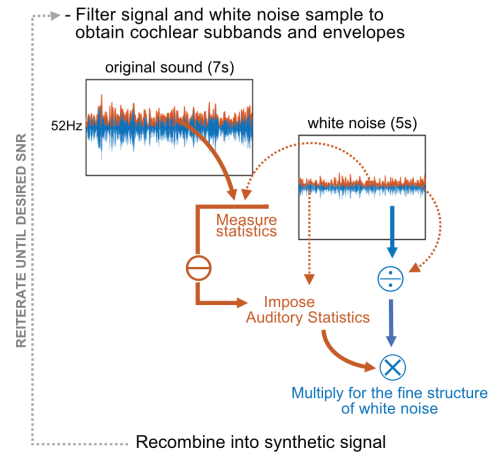
1177

1178 **Figure 2-1. Auditory Evoked response for repeated and novel sounds. Related to**
1179 **Figure 2.** (A) Grand-average topographies across participants of the responses to
1180 standard and oddball sounds for each experiment (Local and Global Discrimination),
1181 displayed for short and long durations (478) at latencies of interest. (B) Grand-average
1182 ERPs across participants of the average amplitude of the central channels displayed in
1183 the legend (red circles on the sketch of a scalp). ERPs are shown for both standard and
1184 oddball sounds for each experiment and duration. Shaded regions show interpolated
1185 standard error of the mean (SE) at each point.
1186

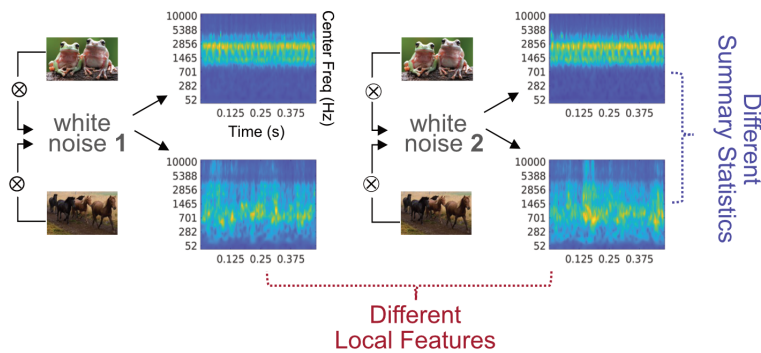
A Computational Texture Model to extract auditory statistics



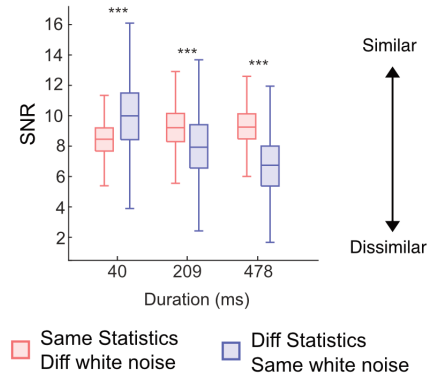
B Schematic of Sound Synthesis



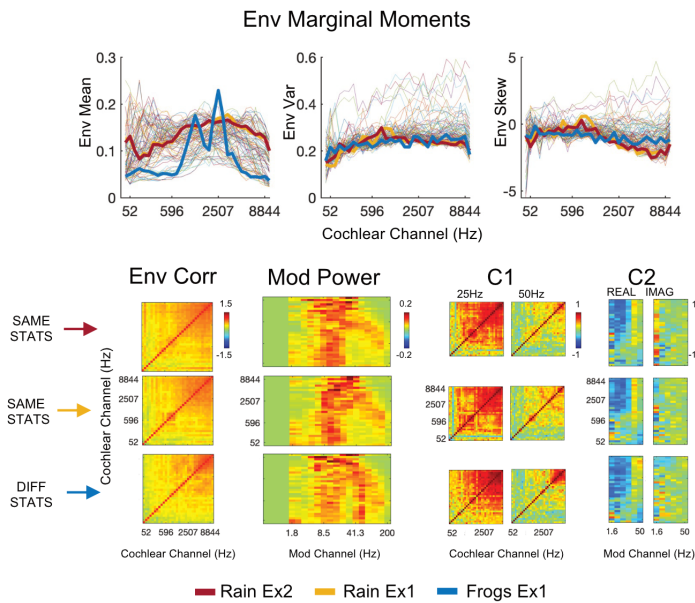
C Impact of white noise and imposed statistics on synthetic sounds



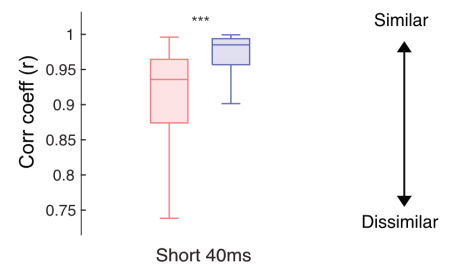
D Similarity of statistics between excerpt pairs



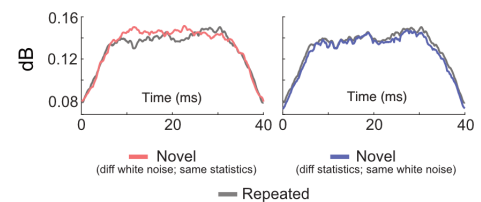
E Comparing auditory statistics of 478ms synthetic sounds



F Similarity of envelopes between short excerpt pairs

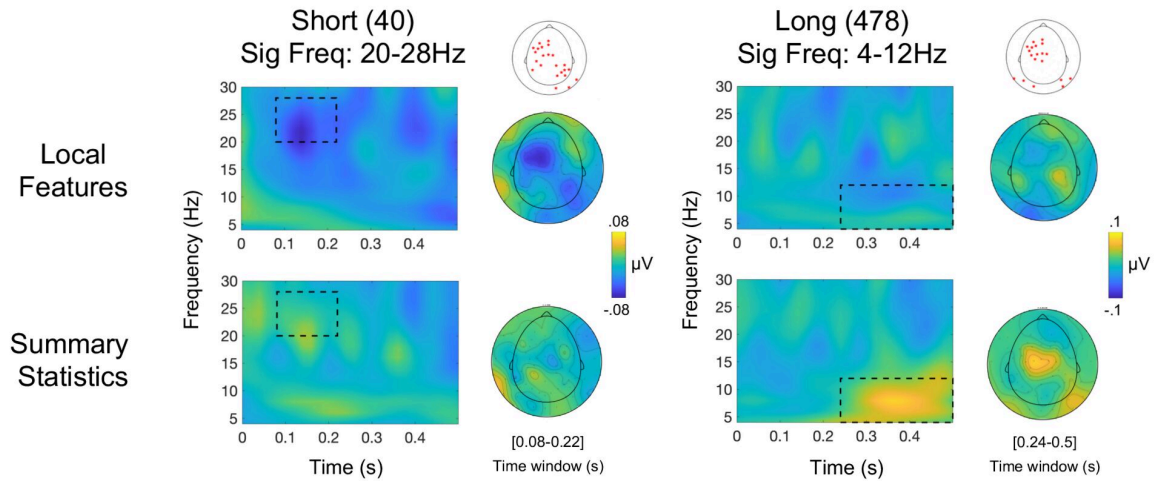


Examples of 40ms envelopes



A

Time-Frequency Results (novel-repeated)



B

Average Power at different Temporal Scales (novel-repeated)

