1
2
3
4
5
6 **A Neural Ensemble Correlation Code for Sound Category Identification**
7
8
9 Authors: Mina Sadeghi[1], Xiu Zhai[1,3], Ian H. Stevenson[2-3], and Monty A. Escabí[1-3]
10
11 Affiliation: [1]Electrical and Computer Engineering, [2]Psychological Sciences, and
12 [3]Biomedical Engineering, University of Connecticut, Storrs, CT 06269
13
14 Correspondence: Monty A. Escabí
15 University of Connecticut
16 Electrical and Computer Engineering
17 371 Fairfield Way U4157
18 Storrs, CT 06269
19
29
30
31
32
33
34
35
36
37

## ABSTRACT

Humans and other animals effortlessly identify sounds and categorize them into behaviorally relevant categories. Yet, the acoustic features and neural transformations that enable the formation of perceptual categories are largely unknown. Here we demonstrate that correlation statistics between frequency-organized cochlear sound channels are reflected in the neural ensemble activity of the auditory midbrain and that such activity, in turn, can contribute to discrimination of perceptual categories. Using multi-channel neural recordings in the auditory midbrain of unanesthetized rabbits, we first demonstrate that neuron ensemble correlations are highly structured in both time and frequency and can be decoded to distinguish sounds. Next, we develop a probabilistic framework for measuring the nonstationary spectro-temporal correlation statistics between frequency organized channels in an auditory model. In a 13-category sound identification task, classification accuracy is consistently high (>80%), improving with sound duration and plateauing at ~ 1-3 seconds, mirroring human performance trends. Nonstationary short-term correlation statistics are more informative about the sound category than the time-average correlation statistics (84% vs. 73% accuracy). When tested independently, the spectral and temporal correlations between the model outputs achieved a similar level of performance and appear to contribute equally. These results outline a plausible neural code in which correlation statistics between neuron ensembles of different frequencies can be read-out to identify and distinguish acoustic categories.

## INTRODUCTION

In the peripheral auditory system, neurons appear to encode sounds by decomposing stimuli into cardinal physical cues, such as sound pressure and frequency, which retain detailed information about the incoming sound waveform. In mid-level auditory structures, such as the inferior colliculus (IC), sounds are then further decomposed into higher-order acoustic features such as temporal and spectral sound modulations. Rather than firing when a frequency is simply present in the sound, neurons in the IC also respond selectively to spectro-temporal structure in the envelopes of frequency channels [1]. In natural sounds, spectro-temporal modulations are highly structured and varied, and the envelopes are correlated both across frequencies and time [2-4]. While low-level cues, such as the sound spectrum, contribute to many auditory tasks, including sound localization and pitch perception [5,6], spectral cues alone are insufficient for identifying most environmental sounds [2]. Manipulating higher-order statistics related to the spectro-temporal modulations of sounds can dramatically influence sound recognition [2]. Temporal modulations contribute to the familiarity and recognition of natural sounds such as running water sounds [7,8], and temporal coherence across frequencies plays a central role in auditory stream segregation [9]. Here we examine to what extent correlated structure in the modulation envelopes of natural sounds is reflected in the correlated structure of neural ensembles and to what extent both neural and sound correlations may be useful for sound category identification. Although correlations are often thought to lead to less efficient representations of individual sensory stimuli [10,11], here we find evidence that correlation statistics may contribute to sound recognition and the formation of acoustic categories.

The spectro-temporal correlation structure of amplitude modulations in sounds is known to contribute to auditory perception [2] and strongly modulates single neuron activity in the auditory midbrain [1]. For single neurons, correlated sound structure can improve signal detection [12], coding of spectro-temporal cues [13] and activates gain control mechanisms [14,15]. However, how correlations in sound modulations are reflected in the correlations in neural ensemble is less

84  clear. Pairs of neurons at multiple levels of the auditory pathway show correlated firing that
85  strongly depends on the spatial proximity of neurons, receptive field similarity, and behavioral
86  state [16-18], but, it remains to be seen is whether these correlations are stimulus-dependent and
87  whether they are detrimental or beneficial for coding. Theories based on efficient coding
88  principles have proposed that stimulus-driven correlations should be minimized in order to
89  minimize redundancies in the neural representation [10,11]. And noise correlations, which reflect
90  coordinated firing in an ensemble of neurons that is not directly related to the sensory stimulus,
91  are thought to directly limit the encoding of sensory information [19,20]. On the other hand,
92  correlations between neurons may be functionally important and have previously been
93  considered as plausible mechanisms for sound localization [21] and pitch identification [22,23]. Here
94  we consider a more general role for stimulus-driven ensemble correlations and test whether they
95  may be useful for category identification and sound recognition.
96          The inferior colliculus receives convergent input from various brainstem nuclei and has
97  the potential to consolidate ascending auditory information to form compact mid-level auditory
98  representations. Here we test the hypothesis that sound correlation statistics alter the correlations
99  between frequency organized neuron ensembles in the IC and that stimulus-driven correlation
100 structure can directly contribute to sound recognition. First, we demonstrate that neural ensemble
101 correlation statistics in the auditory midbrain are strongly affected by sound correlation statistics.
102 We then show that the resulting neural ensemble statistics can be used to discriminate stimuli,
103 and that the sound correlation structure, alone, can be used to identify and group sounds into
104 distinct categories. In analyzing both neural and sound correlations we consider the role of
105 spectro-temporal correlations, generally, as well as the specific contributions of temporal and
106 spectral correlations on sound categorization. We find that stimulus-driven correlations between
107 neuron ensembles in the auditory midbrain may provide a signature for the nervous system to
108 recognize and categorize natural sounds.
109

## RESULTS

111          We first demonstrate that neural correlation statistics between frequency organized
112 (tonotopic) neural ensembles in the auditory midbrain are highly structured containing both time
113 and frequency dependent information that can be used to recognize sounds. Using a neurally-
114 inspired sound representation, we then characterize the modulation correlation statistics of
115 various sound categories and demonstrate how such statistics could be read out by the auditory
116 system for sound category identification.
117

### Decoding Neural Ensemble Correlation Statistics to Recognize Sounds

119          We used multi-channel, multi-unit neural recordings in the auditory midbrain (inferior
120 colliculus, IC) of unanesthetized rabbits to characterize and determine whether neural
121 correlations between recording sites in the IC are affected by the correlation structure of natural
122 sounds and whether such neural ensemble statistics could potentially contribute to sound
123 recognition.
124          The spatio-temporal correlation statistics from an example penetration site demonstrate
125 the diversity of neural ensemble correlation statistics observed in a frequency organized
126 recording site. As expected for the principal nucleus of the IC, frequency responses areas are
127 tonotopically organized varying from low to high frequency with recording depth (Fig. 1A; ~2-
128 12 kHz; low frequencies more dorsal, high frequencies more ventral) and, for this reason, spatial
129 correlations are referred to as spectral correlations in what follows. After cross-correlating the

130  outputs of each recording site (see Materials and Methods) we find that neural ensemble
131  correlation statistics to natural sound textures are highly diverse (Fig. 1H, I). The spectral
132  correlation matrix (at zero-time lag) reflects the instantaneous correlation of the neural activity
133  between different frequency organized recordings sites (Fig. 1C for fire; 1F for water) and is
134  analogous to the frequency coincidence detection previously proposed for pitch perception [22]. In
135  contrast, the autocorrelograms for each recording site reflect the temporal correlation structure of
136  the neural activity at different recording locations (Fig. 1D for Fire; 1G for water). The neural
137  response correlation structure of the same penetration site for five natural sounds (crackling fire,
138  bird chorus, speech babble, running water, rattling snake) are shown both along the spectral (Fig.
139  1H; i.e., different recording sites with different best frequencies) and temporal dimensions. We
140  find distinct patterns for both spatial/spectral (Fig. 1H) and temporal (Fig. 1I) correlations. For
141  instance, in this penetration site, the water and crowd have a similar spatial correlation matrices
142  with the highest correlations localized to neighboring low frequency channels (<4 kHz, along the
143  diagonal). Correlations are much more extensive and widespread for the fire and bird chorus,
144  extending across the entire array. Nearby channels are still highly correlated (along the
145  diagonal), but distant recording channels also have high correlations. The rattling snake sound,
146  by comparison has its own distinct correlation patterns where channels with best frequencies
147  above ~2 kHz are strongly correlated with one another.

148      In the time-domain, the temporal correlations of the neural activity likewise have unique
149  patterns and time-scales for each sound (Fig. 1I). The temporal correlations for the crowd and
150  water sounds are relatively fast (Fig. S3C; 3.0 ms, width at 50% of the maximum; 4.9 and 5.9 ms
151  width at 10% of the maximum). By comparison, although the bird chorus sounds have a similar
152  sharp temporal correlation (4.9 ms half width), there is also a broader and substantially slower
153  component (112 ms, 10% width). The crackling fire produces similarly fast neural correlations
154  (2.9 ms half width; 6.9 ms, 10% width), however these extend across the entire recording array
155  (~1-8 kHz). This pattern likely arises because of the impulsive character of the popping ambers
156  which are both brief in time and broadband in frequency. By comparison, although the rattling
157  snake sound also has a precise temporal correlation at zero lag (3.9 ms, 50% width; 8.9 ms, 10%
158  width) the ensemble generates a broad periodic pattern with a period of ~50 ms reflecting the
159  structure of the rattling at ~20 Hz.

160      Since the structure of the neural ensemble correlations is highly diverse, we next quantify
161  to what extent these statistics could be used to identify sounds (Fig. 2). We use a cross-validated
162  minimum distance classifier (see Materials and Methods) to determine whether the spectral or
163  temporal neural correlation structure could distinguish amongst the five sounds delivered. Neural
164  classifier performance results for spectral, temporal, and spectro-temporal classifiers are shown
165  for the penetration site shown in Fig. 1 as a function of sound duration (Fig. 2A). For short
166  durations (62.5 ms), the spectro-temporal classifier performance for each sound is quite variable,
167  although the average performance is above chance level (20%). As expected, the individual
168  sound and average classifier performance improves with the duration of the recording. When we
169  consider spectral correlations in isolation, the spectral classifier performance is nearly identical
170  to the full classifier (spectro-temporal) and has similar trends for each of the individual sounds.
171  By comparison, the temporal classifier has somewhat poorer performance.

172      Although similar average trends are observed for both the temporal and spectral neural
173  correlation, differences are observed between the performance for different sounds. For example,
174  the fire sound has 100% classification accuracy for the temporal classifier regardless of temporal
175  resolution used while the spectral classifier accuracy is near 0% at 62.5 ms duration and quickly

176  transitions to 100% at ~200 ms duration (Fig. 2A). In general, we observe similar differences
177  across the population where some sounds are better classified with temporal correlations
178  compared to spectral correlations or vice versa (Fig. S1). For example, the average population
179  performance of the temporal correlation classifier for snake remained near chance level
180  regardless of the neural response duration (Fig. S1C) while the spectral correlation classifier
181  performance is comparatively high and improved with increasing response duration (Fig. S1B).
182  Examination of the neural correlations and the acoustic correlations in the rattling snake sound
183  reveals that the periodicity of the snake rattle for the training (first half) and validation (second
184  half) data has different modulation frequencies (~16.5 Hz vs. 20 Hz) which leads to
185  misclassification for a number of penetration sites. By comparison, the spectral correlation
186  structure for this same sound is stable and sufficiently unique to allow high classification
187  performance.
188      While the temporal classifier performance for the example penetration site is relatively
189  high, it is possible that frequency-dependent information contributes to the temporal classifier
190  performance. The temporal correlation matrix used for the classification maintains the tonotopic
191  ordering of the recording channels, and this frequency-dependence could be used to improve
192  classification. Altering tonotopy by manipulating the frequency ordering of sounds, for instance,
193  has been shown to impair pitch perception and vowel recognition [6,24]. To remove the frequency
194  dependence and measure the contribution of purely temporal correlations, we distorted the
195  tonotopic ordering from the temporal correlation matrix. Here we randomize the ordering of the
196  recording channels during the classification. Doing so substantially reduces the classifier
197  performance from nearly 100% average performance to ~50% for the longest sound duration
198  measured (1 s; Fig. 2A). Thus, removing the tonotopic ordering to isolate purely temporal
199  correlations leads to a substantial reduction in the classifier performance.
200      Averaging across all penetration sites (N=13) we find that, just as with the single
201  example, neural correlation structure is highly informative and can be used as a neural response
202  feature to recognize sounds. Across all sounds and classifiers, the average neural classifier
203  performance improves with increasing sound duration (Fig. 2B). As with the example
204  penetration site (Fig. 2A), the spectral and temporal classifier performance can be very different
205  for different sounds (Fig. S1). Distinct differences in the overall classification accuracy and
206  confusions are also observed in the average confusion matrices (Fig. 2C, average across all
207  penetration sites for 1 s duration sound). Furthermore, when the tonotopic organization is
208  removed from the temporal classifier the classifier accuracy is substantially reduced (Fig. 2B, far
209  right). The average performance trends of the spectro-temporal classifier differs substantially
210  from the temporal classifier but followed a nearly identical trend across sounds as the spectral
211  classifier (Fig. S1, compare A, B and C). This indicates that the performance of the spectro-
212  temporal classifier is dominated by the spectral correlations. The performance of the temporal
213  classifier (Fig. 2B), by comparison, is somewhat lower than the spectral classifier and removing
214  the tonotopic organization to isolate strictly temporal correlations further reduces the classifier
215  accuracy (Fig. 2B, far right). These differences suggest that spectral and temporal correlation
216  statistics can contribute differentially to sound identification and that, when both are combined,
217  spectral correlations dominate the classification performance. Overall, we find that spectral and
218  temporal correlations in IC neural ensembles are highly structured and could serve as features for
219  sound recognition.
220
221  **Correlation Statistics of Natural and Man-Made Sound Categories**

222    After establishing that neuron ensembles in IC have highly structured spatio-temporal
223    correlation statistics, we now aim to determine to what extent these correlations are inherited
224    from the correlations in the sounds themselves. Here we use an auditory model to characterize
225    the structure of spectro-temporal correlations in an assortment of natural sound categories and
226    develop a Bayesian classifier to determine their potential contribution towards sound category
227    identification. Natural sounds representing 13 acoustic categories along with white noise (as a
228    reference) are analyzed with the auditory model. We first describe the time-averaged correlations
229    to identify structural differences between sound categories. We then analyze short-term, time-
230    varying correlations to characterize the temporal dynamics and nonstationarity of different sound
231    categories.
232
233    *Average Correlation Statistics and Diversity*
234    Here we evaluate correlations between modulations of different frequency-selective
235    outputs of a cochlear model (see Materials and Methods). The time-averaged spectro-temporal
236    correlations (Fig. 3B and F) highlight distinct acoustic differences between sounds. At zero-lag,
237    the spectro-temporal correlations reflect the instantaneous correlations between different
238    frequency channels (Fig. 3C and G), while correlations within the same frequency channels at
239    different lags reflect the autocorrelations (Fig. 3D and H). For example, the spectral correlation
240    structure of speech is relatively broad reflecting the strong commodulation between frequency
241    channels. This contrasts the running water excerpt, where the correlations are largely
242    diagonalized, indicating minimal correlation between distant frequency channels. The temporal
243    correlation structure in speech exhibits a relatively slow temporal structure (64 ms half width),
244    which reflects the relatively slow time-varying structure of speech elements and words [25,26]. By
245    comparison, the water sound has a relatively fast temporal correlation structure (4 ms half width)
246    that is indicative of substantially faster temporal fluctuations in the sound power.
247    The average spectral (Fig. 4A) and temporal (Fig. 4B) correlations of 13 sound categories
248    and white noise reflect conserved acoustic structure in each of the categories, and, for the
249    categories tested here, this structure is highly diverse. Certain sounds categories, particularly
250    background sounds such as those from water (rain, running water, waves) and wind, have
251    relatively restricted spectral correlations (diagonalized) and relatively fast temporal correlation
252    structure (impulsive; correlation half width: water=3.9 ms; wind=3.6 ms). Such diagonalized and
253    fast temporal structure reflects the fact that these sounds are relatively independent across
254    frequency channels and time. Note that due to the bandwidth of the overlapping filters in the
255    cochleogram, white noise has similar, restricted spectral correlations, rather than perfectly
256    uncorrelated frequency channels. Other sounds, such as isolated vocalizations (e.g., cat, dogs and
257    speech) have more varied and extensive spectral correlation that are indicative of strong coherent
258    fluctuations between frequency channels. Such sounds also have relatively slow temporal
259    correlation structure (correlation half width: speech=64.9 ms; dogs=76.3 ms; cats=117.3 ms),
260    indicating slow dynamics associated with the production of vocalizations.
261    Although the average statistics illustrate differences between categories that could
262    facilitate recognition, each particular sound in a given category can be statistically quite different
263    from the average. This diversity could potentially limit the usefulness of such statistics in
264    recognition. Sound categories in which the correlation statistics exhibit little diversity from
265    sound-to-sound may be easier to identify, using a template-based classifier, for instance, while
266    sounds with large amount of diversity might be expected to be more difficult to identify. For this
267    reason, we developed a category diversity index (CDI) to quantify the diversity in the correlation

268  structure within each of the sound categories (Fig 4C). Of all sounds in our database, white noise
269  has the smallest diversity, which is expected given that white noise is wide-sense stationary. Of
270  the natural sounds, speech has the lowest diversity (CDI=0.12). Although this is unexpected, it
271  likely reflects the fact that the sound segments used in our database consisted of different speech
272  excerpts from a single male speaker. More generally, however, there is no clear distinction or
273  trend between different classes of sounds. For instance, the diversity indices for isolated
274  vocalization categories are quite varied. Bird songs have a relatively high CDI of 0.29 and
275  barking dogs had intermediate values (0.21). Similarly, the diversity of background sounds is
276  quite varied ranging from highly diverse categories such as fire (CDI=0.39) to less diverse
277  categories such as crowd noise (CDI=0.15) and wind (CDI=0.18). Although such trends partly
278  reflect biases in the selection of sounds in the database, they also likely reflect acoustic
279  properties that are unique to each sound category.
280
281  *Short-Term Correlation Statistics and Stationarity*
282      Although the average correlation statistics provide some insights into the large-scale
283  structural differences between sound categories, many sounds, such as vocalizations, exhibit
284  nonstationary structure with complex temporal dynamics and are not well described by time-
285  averaged correlations. Here, to characterize this nonstationary structure, we use time-varying,
286  modified short-term correlations [27]. Computations involving temporally localized and
287  continuously-varying correlations may be more plausible than time-averaged correlations, since
288  lemniscal auditory neurons integrate sounds over restricted integration time windows
289  (approximately 10 ms in the midbrain to 100 ms in cortex [18]).
290      The short-term correlation decomposition of a speech excerpt and a flowing water sound
291  are shown in Fig. 5 (A-D for speech; E-H for water; see supplementary movies, M1-M4). For
292  these examples, we compute the time-varying correlation functions of the cochleograms using a
293  400 ms moving window, and these two examples illustrate the extreme differences that are
294  possible with the time-varying statistics. Speech is highly nonstationary at this time-scale and the
295  correlation structure (spectro-temporal, spectral and temporal correlations) varies considerably
296  from one instant in time to the next. Such nonstationary spectro-temporal correlation structure is,
297  in part, due to the differences between periods of speech and silence. However, even within
298  speech periods the correlation structure can vary and can be quite distinct for each word. These
299  differences likely reflect the range of articulatory mechanisms involved during speech
300  production and the rapidly changing phonemes, formants, and pitch. In sharp contrast to speech,
301  water sounds are relatively stationary. For an example sound, the correlation statistics at each
302  instant in time are relatively consistent and closely resemble the time-averaged correlation (Fig.
303  5B-D and F-H; average shown in gray boundary). No matter which time segment we examine,
304  spectral correlations are diagonalized indicating that only neighboring channels have similar
305  envelopes, while temporal correlations are relatively fast with similar fast time constants across
306  all frequencies.
307      To quantify the degree of stationarity (or lack of) in the short-term correlation function
308  we developed a stationarity index (SI, see Materials and Methods). The SI quantifies the degree
309  to which the sound correlation function varies dynamically over time with a value of 1 indicating
310  perfect stationarity and a value of 0 indicating that the sound is highly nonstationary. Although
311  the results are quite varied, we note several trends. First, except for fire sounds, environmental
312  sounds (including running water, waves, thunder and wind) tend to have the highest average SI
313  (SI=$0.49 \pm 0.04$; mean$\pm$SEM). This might be expected given that such environmental sounds

314  typically consist of mixtures of randomly arriving sound elements (e.g., water droplets, air
315  bubbles, etc.) and have been shown to be perceptually well described by average statistics [2,7]. By
316  comparison, vocalized sounds have a somewhat lower nonstationary index (SI=$0.28 \pm 0.04$;
317  mean$\pm$SEM) at the analysis time-scales employed. This is expected given that vocalizations
318  transition from periods of silence and vocalizations at time-scales of just a few Hz [26] and even
319  within vocalization segments the correlation structure can vary dynamically from moment-to-
320  moment (e.g., Fig. 5). In the classification analysis that follows, we aim to describe this
321  nonstationary structure and determine to what extent ignoring nonstationary impairs sound
322  categorization.
323

## Decoding Correlation Statistics to Categorize Sounds

325  Given that the neural ensemble correlation statistics are strongly modulated by the
326  correlation structure in sounds and that natural sound ensembles have highly varied spectro-
327  temporal correlation statistics, we next tested whether these high-order correlation statistics
328  could directly contribute to sound category identification. Here we aim to quantify the specific
329  contributions of spectral and temporal correlations as well as the overall accuracy of
330  classification when both features are used. While both spectral and temporal envelope cues can
331  contribute to a variety of perceptual phenomena they often do so differently. For instance,
332  speech recognition can be performed with low spectral resolution so long as detailed temporal
333  information is preserved in each frequency channel [28]. By comparison, music perception requires
334  much finer spectral envelope resolution [29]. Thus, it is plausible that one of the two dimensions
335  (temporal or spectral) could be more informative for specific sound categories or for the sound
336  category identification task as a whole.
337  Rather than using a minimum distance classifier on the time-averaged correlations (i.e.,
338  Fig. 3 and 4), such as the one we used for classifying neural correlations, here we develop a
339  probabilistic model of the short-term correlation statistics. After computing the time-varying
340  correlations for purely spectral, temporal, or spectro-temporal correlations (Fig. 5), we reduce the
341  dimensionality of the features using principal component analysis (Fig. S4). We then fit the low-
342  dimensional representation of the correlations using an axis-aligned Gaussian mixture model for
343  each sound category (see Materials and Methods). After training, we classify test sounds by
344  comparing the posterior probability of the sounds under the mixture models of each category. By
345  limiting the short-term correlations to spectral (Fig. 5C and G), temporal (Fig. 5D and H), or
346  spectro-temporal (Fig. 5B and F) we can measure how each of these acoustic dimensions
347  contributes to categorizing sounds (see Materials and Methods). Ignoring the nonstationary
348  structure by averaging the statistics over time impaired the sound category performance,
349  particularly for short duration sounds where the nonstationary classifiers had ~25% higher
350  accuracy compared to the classifier based on averages (Fig. S5). The time-varying statistical
351  structure of the sounds can, thus, contribute to more accurate sound categorization.
352  We optimized the model and classifier for each task separately using multiple temporal
353  resolutions ($\tau_W$ =25-566 ms; Fig. 6A) while using the maximum sound duration (10 s). The
354  optimal resolution for both the temporal and spectral correlation classifiers is 141 ms, while the
355  optimal resolution for the joint spectro-temporal correlation classifier is slightly faster (100 ms).
356  For both spectral, temporal and spectro-temporal correlations, the classifier performance is above
357  chance for all temporal resolutions tested (Fig. 6A; chance performance = 7.69%; $p<0.01$, t-test,
358  Bonferroni correction). For the spectro-temporal classifier, performance improves with
359  increasing sound duration reaching a maximum accuracy of 84% (Fig. 6B). Individually, spectral

360 and temporal correlations alone achieve similar maximum performance (spectral=83%,
361 temporal=81%; p=0.60, two-sample t-test) and neither is significantly different from the overall
362 spectro-temporal performance (spectral versus spectro-temporal: p=0.79, two-sample t-test;
363 temporal versus spectro-temporal: p=0.42, two-sample t-test). This indicates that that both
364 spectral and temporal correlations contribute roughly equally to the sound identification task for
365 the full sound duration. However, the spectral classifier performance increases at a faster rate
366 than the temporal classifier (reaching 90% of maximum in 1.7 vs. 3.0 s; Fig. 6B), indicating that
367 evidence about the sound category may be accumulated more efficiently using spectral
368 correlations. The joint spectro-temporal classifier improves at an even faster rate (reaching 90%
369 of its maximum in 1.2 s).
370 As with the neural correlations, we find distinct differences between the performance of
371 the spectral and temporal classifiers for the individual sound categories (Fig. 6C and D). The
372 recognition accuracy for each sound and the types of confusions that occur are quite different for
373 the spectral and temporal classifiers (Fig 6C). Certain sounds such as waves perform
374 substantially better for the spectral correlation classifier (spectral=100% vs temporal=33%).
375 Other sounds, such as wind exhibit higher performance with the temporal classifier
376 (spectral=47% vs. temporal=80%). Thus, although on average the performance of the spectral
377 and temporal correlations classifier is comparable, performance of certain sound categories
378 appears to be dominated by one of the two sets of features (Fig 6D).
379 We also evaluate the performance of the classifier in a two-alternative forced choice task
380 where we require the classifier to distinguish vocalization and background sound categories. The
381 model performance is consistently high with accuracy rates for the spectro-temporal classifier
382 approaching 90% for background sounds and nearly 100% for vocalizations (Fig. 7).
383 Interestingly, the performance for identifying vocalizations improves with increasing sound
384 duration while the performance for background sounds remains constant. This behavior occurs
385 regardless of whether temporal, spectral or spectro-temporal correlations are used (Fig. 7, A-C).
386 Since background sounds are more stationary, their statistics can be assessed quickly in this
387 relatively simple task by the classifier. Vocalizations, on the other hand, are nonstationary over
388 longer time-scales and have epochs of silence that may require the classifier to accumulate
389 evidence over longer time. Additionally, the classification accuracy of background sounds is
390 comparable (~90%) for spectral and temporal features. Except for the shortest sound duration
391 used (100 ms), vocalization sounds are more accurately classified with spectral correlations
392 compared to temporal correlations (p<0.01, t-test with Bonferroni correction), and the spectral
393 correlations alone appear to account for most of the performance of the spectro-temporal
394 classifier.

## DISCUSSION

397 Here we have demonstrated that natural sounds can have highly structured spectro-
398 temporal correlations and that this acoustic structure induces highly structured correlations
399 between neural ensembles in the auditory midbrain. Stimulus-driven correlation structure, in
400 both time and frequency (or space), is highly informative and conveys information about the
401 sound identity, with evidence accumulation times on the order of a few seconds. Spectral and
402 temporal correlations both contain information about sound categories and, nonstationary
403 correlation structure conveys information beyond the corresponding time-average sound
404 correlations. Altogether, these results demonstrate how time-varying correlations between

405 modulations of sounds and the resulting correlations in frequency organized neural ensembles
406 have the capacity to contribute to sound recognition.
407
408 *Using Sound Driven Correlation Statistics for Recognition*
409       Previous work on how populations of neurons encode stimuli has emphasized the distinct
410 roles of noise-driven and stimulus-driven correlations between neurons in limiting the
411 information capacity of neural ensembles. Noise correlations, which are the result of coordinated
412 firing not directly related to the stimulus structure, are thought to limit the encoding of sensory
413 information [19,20]. At the same time, theories of efficient coding suggest that neural ensembles
414 should minimize the amount of stimulus-driven correlated firing to limit redundancies in the
415 neural code [10,11,18]. Here we find that stimulus-driven correlations between neural ensembles in
416 the IC are highly structured. In contrast, the noise correlations in our neural data are largely
417 unstructured. They are localized primarily to nearby recording channels as well as to brief time-
418 epochs, and they do not vary systematically with the sound (Fig S2.1-2.3 and S3). These finding
419 thus provide an alternative and, perhaps, underappreciated viewpoint: correlations in neuron
420 ensembles can convey critical information about the stimulus itself. Rather than redundant
421 structure that should be removed [11], sound correlation statistics can be viewed as high-level
422 acoustic features or cues that are highly informative and which, as demonstrated, drive correlated
423 activity in mid-level auditory structures.
424       Indeed, recent studies on sound textures confirm that correlation structure is a critical cue
425 required to create realistic impressions of sounds [2]. And there is growing evidence from neural
426 recordings that neural correlations can convey information about stimuli or behavior during
427 decision making [30-32]. Our findings extend these views, by demonstrating that neural ensembles
428 in a mid-level auditory structure can directly signal sound correlation structure and that sound
429 categories have unique correlation statistics that may promote sound categorization.
430
431 *Biological Plausibility*
432       While our results demonstrate that neural ensemble correlations in the auditory midbrain
433 are highly structured in both the spectral and temporal domains, how and if such information is
434 used by higher brain regions needs further exploration. One possible mechanism for representing
435 correlations between neurons has been previously considered for pitch detection: the frequency
436 coincidence detection network [22,23]. The key proposal of this network is that neurons encoding
437 different frequencies project onto the same downstream neurons that then detect coincident
438 frequencies. Given that anatomical connections within and beyond the IC can span a broad range
439 of frequencies, the central auditory system anatomy allows for such a possibility. Many
440 projection neurons in the IC have collaterals that extend across multiple frequency-band lamina
441 and send their outputs to auditory thalamus [33]. Connections between thalamus and auditory
442 cortex as well as intra-cortical connections, although frequency specific, can also extend across
443 several octaves [34,35]. This pattern of connectivity has the capacity to integrate information across
444 frequency channels and may subserve a coincidence-like operation between inputs of different
445 frequency. This coincidence detection would allow downstream neurons to directly compute
446 spectral correlations, not just for pitch detection but, for sound categorization as well.
447
448 *Contribution of Spectral and Temporal Correlations*
449       A key question brought up the model and neural data is the extent to which spectral and
450 temporal correlations individually contribute to sound recognition. The auditory model suggests

451　that both can contribute roughly equally for the sound categories tested, and performance is only
452　marginally better when the two dimensions are combined. On the other hand, although the neural
453　classification performance favors spectral (i.e. spatial) correlations over temporal, both sets of
454　features have discriminative power. However, when the tonotpic organization is removed so that
455　purely temporal correlations are used for the neural data, the temporal classifier performance is
456　substantially reduced. This reduced performance is not observed in the model, since the temporal
457　classifier used is specifically designed to isolate strictly temporal cues (see Materials and
458　Methods). The observed difference between model and neural data may partly reflect differences
459　in the sound paradigm and classifier used (single sound identification versus category
460　identification). It is also possible that such differences are partly attributed to neural mechanisms
461　not captured by the model, including spectro-temporal nonlinearities [1] or adaptation [36] at the
462　level of the IC. As such, more extensive studies are necessary to parse out differences between
463　the two domains.
464　　　　Insofar as neural mechanisms for computing the sound correlation statistics, spectral
465　models have broader biological support. Spectral correlations could be computed using simple
466　coincidence detection, where two sound modulated inputs tuned to distinct frequencies converge,
467　multiplicatively on downstream neurons. This type of spectral convergence is widespread in
468　auditory system anatomy and the required multiplicative interactions have been previously
469　described [37,38]. On the other hand, computing temporal correlations requires coincidence
470　detection of activity occurring at different times. Delay lines, differences in integration times, or
471　feedback loops could all subserve these computations (with increasing temporal differences), but
472　these phenomena are more speculative and lack strong anatomical or physiologic evidence.
473
474　*Resolution and Integration Time-Scales for Feature Analysis and Inference*
475　　　　Sound categorization performance for both the neural data and cochleogram model
476　correlations improve over the course of 1-2s and depended strongly on sound duration, similar to
477　human listeners [39]. This brings up the question of whether the previously observed perceptual
478　integration times in human observers should be attributed to a central neural integrator which
479　averages sound statistics with this relatively long time-constant of a few seconds. Rather than
480　computing average statistics about a sound over long time-scales, its plausible that sound
481　statistics themselves are integrated and estimated at relatively short time scales analogous to the
482　optimal integration resolution of our model (~150 ms). The long time-scales of a few seconds
483　required to make perceptual decisions, may instead reflect a statistical evidence accumulation
484　process, as previously proposed for cortical areas involved in decision making [40].
485　　　　In modeling correlations in natural sounds, we find that the times required to accumulate
486　statistical information about sound categories are roughly an order of magnitude larger than the
487　optimal temporal resolution for calculating the correlations (~100-150 ms). This is consistent
488　with a temporal resolution-integration paradox previously observed for neural discrimination of
489　sounds [41]. However, these time-scales are substantially longer than those previously reported for
490　neural discrimination of sounds in auditory cortex. Previous studies identified an optimal
491　temporal resolution of ~10 ms and an integration times of ~500 ms for discriminating pairs of
492　sounds in auditory cortex [41,42]. It is likely that such differences in temporal resolution and
493　integration times largely reflect differences in the encoded features for our categorization
494　paradigm and differences in how statistical evidence for the task is accumulated. The timescales
495　for auditory cortex in these previous studies were optimized for the discrimination of pairs of
496　sounds based on spike train distance measures at a single neuron level. By comparison, here we

497    use high-order correlation statistics of either frequency-tuned cochlear channels or neuron
498    ensembles as the primary feature for categorizing many more sounds.
499
500    *Conclusion*
501          We have shown how the spectro-temporal correlation structure of natural sounds shows
502    reliable differences between categories. For nonstationary sounds, this structure fluctuates on
503    relatively fast (~150 ms) timescales even within a single sound. Surprisingly, this acoustic
504    structure is reflected in the correlation structure of neural ensembles and can be used for accurate
505    categorization with evidence accumulation time-scales of a few seconds. Together, our results
506    from neural data and sound statistics suggest that spectro-temporal correlations in the auditory
507    system may play an important role, not just in pitch perception or localization, but in sound
508    recognition more generally.
509

510 Figure 1: Neural ensemble correlation statistics for an auditory midbrain penetration site. (A)
511 Neural recording probe and the corresponding frequency response areas at 8 staggered recording
512 sites show tonotopic organization (red indicates high activity, blue indicates low activity). aMUA
513 activity for the 16-recording channels for a (B) fire and (E) water sound segment (red indicates
514 strong response, blue indicates weak response). The spectral (C=fire; F=water) and temporal
515 (D=fire; G=water) neural ensemble correlation for the penetration site. (H) Spectral and (I)
516 temporal correlations of the recording ensemble show distinct differences and unique patters
517 across the five sounds tested.

518

519 Figure 2: Using neural ensemble correlation statistics to identify sounds. (A) Classification
520 results for the penetration site shown in Fig. 1. The average classifier performance (red curve)
521 and the performance for each individual sound (gray lines) are shown as a function of the sound
522 duration for four classifiers. In all cases, classifier performance improves with sound duration.
523 The combined spectro-temporal classifier has the highest performance, followed by the spectral
524 and temporal classifier. Removing the tonotopic ordering of recording sites for the temporal
525 classifier (far right) substantially reduces its performance. (B) Average performance across all of
526 the IC penetration sites shown as a function of sound duration. Gray bands represent SD. (C)
527 Average confusion matrices obtained across all penetration sites for the corresponding conditions
528 shown in (B) for a sound duration of 1 second.

529

530 Figure 3: Measuring the average correlation structure of natural sounds. Illustrated for a speech
531 and a flowing water sound. The spectro-temporal correlations are obtained by cross-correlating
532 the frequency organized outputs of a cochlear model representation (A, E). The resulting spectro-
533 temporal correlation matrices (B, F) characterize the correlations between frequency channels at
534 different time-lags. The spectro-temporal correlations are then decomposed into purely spectral
535 (C, G) or temporal (D, H) correlations. Speech is substantially more correlated across frequency
536 channels and its temporal correlation structure is substantially slower than for the water sound.

537

538 Figure 4: Sound correlation statistics for the thirteen sound categories and white noise. The
539 category average (A) spectral correlation matrix and (B) temporal correlations show unique
540 differences amongst the thirteen sounds examined. (C) The category diversity index (CDI)
541 quantifies the variability of the correlation statistics for each category. A CDI of 1 indicates that
542 the sound category is diverse (the correlation statistics are highly variable between sounds) while
543 0 indicates that the category is homogenous (all sounds have identical correlation statistics).

544

545 Figure 5. Short-term correlation statistics and stationarity. The short-term correlation statistics
546 are estimated by computing the spectro-temporal correlation matrix using a moving sliding
547 window. The procedure is shown for an excerpt of (A) speech and (B) water (additional
548 examples in supplementary movies, M1-M4). The sliding window (400 ms for these example) is
549 varied continuously over all time points but is shown for three select time points for this
550 example. The short-term statistics are also shown for the spectral and temporal correlation
551 decompositions. Note that for speech, the correlations change dynamically from moment-to-
552 moment and differ from the time-average correlations (gray panel) indicating nonstationary
553 structure. By comparison, the time-varying correlations for water resemble the time-average
554 correlations (gray panel) indicating more stationarity. (C) Stationarity indices for the thirteen

555  categories and white noise. Speech has lowest stationarity values while white noise is the most
556  stationary sound.
557
558  Figure 6: Using short-term correlation statistics to categorize sounds in a 13-category
559  identification task. A cross-validated Bayesian classifier is applied to the short-term correlation
560  statistics (spectral, temporal and spectro-temporal) to identify the category of each of the test
561  sounds (see Materials and Methods). (A) Both the spectral and temporal classifiers had an
562  optimal temporal resolution of 144 ms (i.e., short-term analysis window size). The optimal
563  resolution of the spectro-temporal classifier, by comparison is slightly higher (100 ms). (B) For
564  all three classifiers, the performance improves with the sound duration. The spectro-temporal
565  classifier performance improved with sound duration at the fastest rate, while temporal
566  correlations had the slowest rate of improvement. (C) Confusion matrices for the three classifiers
567  for 10 s sound durations. (D) Performance for the three classifiers shown as a function of sound
568  category (measured at the optimal resolution and at 10 s sound duration).
569
570  Figure 7: Classification performance in a two-category identification task distinguishing
571  vocalization from background sound categories. For all three classifiers, the overall performance
572  is consistently high and improved with increasing sound duration. Vocalization classification
573  performance is highest for the (C) spectro-temporal classifier and shows a nearly identical trend
574  for the (A) spectral classifier. The performance of the temporal classifier, however, is 20 %
575  lower. For background sounds, classification performance did not improve over time and is
576  consistently high (~90%) for all three classifiers.
577
578
579
580
581
582
583
584
585

## MATERIALS AND METHODS

*Animal Experimental Procedures*

Animals are handled according to approved procedures by the University of Connecticut Animal Care and Use Committee and in accordance with National Institutes of Health and the American Veterinary Medical Association guidelines. Multi-channel neural recordings are performed in the auditory midbrain (inferior colliculus) of unanesthetized female Dutch-Belted rabbits ($N = 2$; 1.5-2.5 kg). Rabbits are chosen for these experiments since their hearing range is comparable to that of humans, and they sit still for extended periods of time which enables us to record from different brain locations daily over a period of several months.

*Surgery*

All surgical procedures are performed using aseptic techniques. Surgical procedures are carried out in two phases with a recovery and acclimation period between procedures. For both procedures, rabbits are initially sedated with acepromazine and a surgical state of anesthesia achieved via delivery of isoflurane (1-3%) and oxygen (~2 liters/min). In the first procedure, the skin and muscle overlying the dorsal surface of the skull are retracted exposing the sagittal suture between bregma and posterior to lambda. Stainless steel screws (0-80) and dental acrylic are used to affixed a brass restraint bar oriented rostro-caudally and to the left of the sagittal suture. Dental acrylic is then used to form a dam on the exposed skull on the right hemisphere between lambda and the interparietal bone. Next, custom fitted earmolds are fabricated for each ear. A small cotton ball is inserted to block the external auditory meatus, and a medical grade polyelastomeric impression compound poured into the ear canal. After ~5 minutes, the hardened impression compound is removed. The ear impression mold is subsequently used to build a cast from which custom ear molds are fabricated.

Following the first surgery and a five-day recovery period, the animal is acclimated over a period of 1-2 weeks to sit still with the head restrained. During this period, the animal is also gradually exposed to sounds through the custom fitter ear molds. Once the animal is capable of sitting still during sound delivery, the second surgical procedure is performed. The animal is again anesthetized as described above, and an opening (~4 x 4 mm) is made on the right hemisphere within the dental acrylic dam and centered approximately 12-13 mm posterior to lambda. At this point, the exposed brain area is sterilized and medical grade polyelastomeric compound is poured into the acrylic dam to seal the exposed region.

*Sound Delivery and Calibration*

Sounds are delivered to both ears via dynamic speakers (Beyer Dynamic DT 770 drivers) in a custom housing and custom fitted ear molds obtained as described above. The molds are fitted with a sound delivery tube (2.75 mm inner diameter) that is connected to the dynamic speaker housing forming a closed audio system. Calibration consisted of delivering a 10-sec long chirp signal at 98kHz sampling rate via TDT RX6 and measuring the audio signal with a B&K calibration microphone and probe tube placed ~5 mm from the tympanum. The measured signal is used to derive the sound system impulse response (via Wiener filter approach; combined speaker driver and tube) and an inverse filter finite impulse response is then derived. Subsequently, all sounds delivered to the animal are passed via the inverse filter which is implemented in real time using a TDT RX6 at 98kHz sampling rate. The sound delivery system has a flat transfer function and linear phase between 0.1 – 30kHz (flat to within ~ ±3 dB).

At each penetration site, we first delivered a pseudo random sequences of tone pips (50

632 ms duration, 5 ms cosine-squared ramp, 300 ms inter-tone interval) covering multiple
633 frequencies (0.1-16 kHz) and sound pressure levels (5-85 dB SPL). These tone-pip sequences are
634 used to measure frequency response areas which allow us to estimate the frequency selectivity of
635 each recording site (different channels).
636        Next, we delivered a sequence of five environmental noises with distinct structural
637 properties to determine whether the ensemble activity of the auditory midbrain reflected the
638 sound frequency dependent correlation statistics present in the sounds. Each sound is 3 s duration
639 and is delivered in a block randomized fashion with a 100 ms inter-stimulus interval between
640 sounds. To avoid broadband transients, the sounds have b-spline onset and offset ramps (20 ms
641 rise-decay time) and are delivered at an RMS sound pressure level of (70 dB SPL). The sounds
642 included: running water, a crackling fire, speech babble, a bird chorus, and a rattling snake
643 sound. These sounds each contain unique time-frequency correlation statistics allowing us to test
644 and quantify whether such statistics are potentially encoded by the auditory midbrain and
645 ultimately represented in the neural ensemble activity. For instance, the water sounds have
646 minimal across-frequency channel correlation since the air bubbles and droplets responsible for
647 this sound are relatively narrow band and occur randomly in time, thus activating frequency
648 channels independently [7]. Sounds such as crackling fire, by comparison, have strong frequency
649 dependent correlations due to crackling embers which produced brief impulsive "pops" that span
650 multiple frequency channels simultaneously. The temporal correlations of these five sounds are
651 also quite varied. For instance, the water sound has a very brief impulsive correlation structure
652 lasting just a few milliseconds whereas the bird chorus and speech babble have a broader and
653 slower temporal correlation function. The rattling snake sound, by comparison, had strong
654 periodic correlations at ~20 Hz.
655
656 *Electrophysiology*
657        Sixteen channel acute neural recording silicon probes (Neuronexus 10 mm probe; 16-
658 linear spaced *recording* sites with 100 um separation; site impedance ~1-3 $M\Omega$) are used to
659 record neural activity from the inferior colliculus of unanesthetized rabbits. We recorded neural
660 data from N=13 *penetration* sites in the inferior colliculus of two rabbits (N=4 and N=9). Since
661 there are 16 recording channels for each penetration site, data is obtained from a total of
662 16x13=208 recording sites within IC. Prior to recording, the polyelastomeric compound is
663 removed from the craniotomy and lidocaine is applied topically to the exposed cortical tissue.
664 The area is then flushed with sterile saline and the acute recording probe inserted at ~12-14 mm
665 posterior to lambda. If necessary, a sterile hypodermic needle is used to nick the dura to allow
666 the electrode to penetrate the neural tissue. An LS6000 microdrive (Burleigh EXFO) is used to
667 insert the neural probe to a depth of ~7.5-9.5 mm relative to the cortical surface where, at this
668 penetration depth, most or all of the recording electrodes are situated in the IC and had clear
669 responses to brief bursts of broadband noise or tones. Neural activity is acquired continuously at
670 sampling rate of 12 kHz using a PZ2 preamplifier and RZ2 real time processor (TDT, Alchua,
671 Florida).
672        The sampled extracellularly recorded neural signals are analyzed offline using an analog
673 representation of multi-unit activity (analog multi-unit activity, aMUA)[43]. From each of the
674 recorded neural traces aMUA is measured by first extracting the envelope of the recorded
675 voltage signal within the prominent frequency band occupied by action potentials spanning
676 frequencies 325 and 3000 Hz (b-spline filter, 125 Hz transition width, 60 dB attenuation) [44]. The
677 bandpass filtered voltage signal is next full-wave rectified and low-pass filtered with a cutoff

678    frequency of 475 Hz (b-spline filter, transition width of 125 Hz and stopband attenuation of 60
679    dB), since neurons in the auditory midbrain typically don't phase lock to envelopes beyond ~500
680    Hz [45]. The resulting envelope signal is next down-sampled to 2 kHz. Such neural envelope
681    signals captures the synchronized activity and the changing dynamics of the local neural
682    population with each recording array in both time and frequency domains [43]. For each recording
683    channel an analog raster is generated which consists of the aMUA response over time and across
684    trials. Each recording had 5 sound conditions where each sound had at least 18 trials (range 18 to
685    39) with a 3 s duration for each trial.
686
687    *Neural Ensemble Stimulus-Driven Correlation Matrix*
688         For each recording penetration, we estimat the *stimulus-driven* correlations of the neural
689    ensemble across the 16-recording channels directly from the measured aMUA signals. The
690    procedure consists of a modified windowed short-term correlation [27] analysis between recording
691    sites in which the correlations are "shuffled" across response trials [18,46]. The windowed
692    correlation approach allows us to localize the correlation function in time, whereas the shuffling
693    procedure is used to remove neural variability or noise from the correlation measurements. Thus,
694    the proposed shuffled windowed correlation allows us to isolate stimulus-driven correlation
695    between recording sites independently of noise driven correlations. The windowed shuffled
696    cross-correlation between the *k*-th and *l*-th recording site is computed as the mean pairwise
697    correlations between different trials of the aMUA envelope:
698

$$\Phi_{kl}(t, \tau) = \frac{1}{N(N-1)} \sum_{m=1}^{N} \sum_{n \neq m} \phi_{mn}(t, \tau)$$

699
700    where *N* is the number of trials in each recording channel (*k* and *l*). $\phi_{mn}(t, \tau)$ is the windowed
701    cross-correlation between the *m*-th and *n*-th response trials in the two channels respectively,
702    where *t* is time and $\tau$ is the cross-correlation delay
703

$$\phi_{mn}(t, \tau) = \int_{-\infty}^{+\infty} r_{k,m}(\gamma) r_{l,n}(\gamma - \tau) W^2(t - \gamma) d\gamma$$

704
705    where $\gamma$ is the time integration variable. Here $r_{k,m}(t)$ is the *m*-th response trial (mean removed)
706    from channel *k* and $r_{l,n}(t)$ is the *n*-th response trial (mean removed) from channel *l*. $W(\gamma)$ is unit
707    amplitude square window centered about $\gamma=0$ of duration *T* sec (range = 62.5 – 1000 ms) which
708    is used to localize the measured signal correlations around the vicinity of the designated time-
709    point (*t*). Note that in this formulation, correlations between recording channels (*k* and *l*) are
710    computed for different response trials (*m* and *n*). The above is implemented using a fast-shuffled
711    correlation algorithm according to [46]:
712

$$\Phi_{kl}(t, \tau) = \frac{1}{N(N-1)} \left( N^2 \cdot \Phi_{PSTH_{kl}}(t, \tau) - \sum_{m=1}^{N} \phi_{mm}(t, \tau) \right)$$

713
714    where $\Phi_{PSTH_{kl}}(t, \tau)$ is the windowed cross-correlation function for the post-stimulus time
715    histograms (PSTHs) between channel *k* and *l*:

716

$$\Phi_{PSTH_{kl}}(t, \tau) = \int_{-\infty}^{+\infty} PSTH_k(\gamma) PSTH_l(\gamma - \tau) W^2(t - \gamma) d\gamma$$

717

718    where the PSTH for the $k$-th and $l$-th channels are

719

$$PSTH_k(t) = \frac{1}{N} \sum_{m=1}^{N} r_{k,m}(t)$$

720

$$PSTH_l(t) = \frac{1}{N} \sum_{n=1}^{N} r_{l,n}(t)$$

721

722    As previously shown [46], this fast-shuffled correlation algorithm resulted in a marked reduction in
723    the computational time ($N + 1$ correlations compared with $N(N - 1)$) for each pair of recording
724    channels. This speedup in the computational time is necessary to bootstrap the data during the
725    model generation and validation applied subsequently (See *Neural Ensemble Classifier*).
726            To remove the influence of the response power on the correlation measurements, the
727    short-term correlation is normalized as a correlation coefficient. This requires that we measured
728    the localized short-term signal variance at each time and delay sample according to

729

$$\sigma_k^2(t) = var[PSTH_k(t)W(t - \gamma)]$$

730

$$\sigma_l^2(t, \tau) = var[PSTH_l(t - \tau)W(t - \gamma)]$$

731

732    The channel-covariance is then obtained as

733

$$c_{kl}(t, \tau) = \frac{\Phi_{kl}(t, \tau)}{\sqrt{\sigma_k^2(t) \cdot \sigma_l^2(t, \tau)}}$$

734

735    Like a correlation coefficient, this population ensemble correlation function is bounded between
736    -1 and 1.

737

738    *Neural Ensemble Classifier*
739            We use a cross-validated neural ensemble classifier to assess whether spatial, temporal,
740    and/or spatio-temporal correlations of the neural ensemble in IC could be used to recognize
741    sounds. Although technically, the correlations are measured between the neural activity across
742    spatially separated electrode channels ($k$ vs. $l$), the electrode channels with our recording
743    paradigm are frequency ordered (Fig. 1A). The neural ensemble activity thus reflects spectral
744    correlations between frequency channels, and, in what follows, we describe spatial correlations
745    between recording channels as spectral correlations.
746            The neural classifier is implemented using a cross validation approach in which half of
747    the data is used for model generation and the second half for model validation. For each
748    recording site, we chose the first half of the aMUA raster for a given sound (1.5 s) as the model

749 and the second half as the validation data. The first and second half of the data are then swapped
750 and the procedure repeated using the first half of the data for validation and the second half for
751 model generation. The model consisted of the time average correlation function at zero-lag for
752 each of the five sounds $c_{kl,s}(0)$, where $s$ indicates the sound. Note that, since the model
753 correlations are averaged across time, $t$ is no longer a variable in the model correlation. The
754 model classification performance is tested and validated by iteratively implementing a minimum
755 Euclidean distance classifier across different validation data segments according to:
756

$$S = \underset{s=1\dots5}{\operatorname{argmin}} \|c_s - c(t)\|$$

757

758 where $s=1\dots5$ are the five sounds tested, S is the classified sound, $c(t)$ is a vector of features
759 obtained from the spectro-temporal correlation function for a particular time segment of the
760 validation data (at time $t$) and $c_s$ is a vector containing the average correlation feature vector for
761 sound $s$ from the model generation data. The center location of each correlation measurement, $t$,
762 is randomly varied by randomly sampling 100 distinct time points ($t$). The analysis is repeated
763 for window sizes ranging between 62.5 and 1000 ms in ½ octave steps. The overall classifier
764 performance is the average across all five sounds and across the 100 randomly selected
765 validation segments.
766     To evaluate the contribution of temporal and spectral correlations for neural classification
767 performance, we implement the neural classifier either using purely temporal, purely spectral, or
768 joint spectro-temporal correlations. The purely spectral classifier only considers correlations at
769 zero lag, $c_{spec}(t) = c_{kl}(t, 0)$, as the primary features (no time lag between different frequency
770 channels). Note that $c_{kl}(t, 0)$ contains strictly frequency dependent information, since the
771 recorded neural channels are tonotopically ordered and delays are removed. For the temporal
772 classifier, we consider the correlations along the diagonal, $c_{temp}(t) = c_{kk}(t, \tau)$, for delays
773 extending between $\tau = -100$ to 100 ms as the primary features. Next, we combined the spectral
774 and temporal correlations to implemented the joint spectro-temporal classifier as follows:
775

$$S = \underset{s=1\dots5}{\operatorname{argmin}} \frac{\|c_{spec,s} - c_{spec}(t)\|}{N_{spec}} + \frac{\|c_{temp,s} - c_{temp}(t)\|}{N_{temp}}$$

776

777 where $N_{spec}$ and $N_{temp}$ are the number of elements contained in the spectral and temporal
778 feature vectors. By normalizing the spectral and temporal distance measures by the number of
779 elements in each vector we ensure that neither the spectral and temporal correlations dominate
780 the categorization.
781     Finally, we implement a purely temporal correlation classifier that lacks frequency
782 organization. Note that for the temporal correlation classifier, the features consist of the envelope
783 autocorrelations taken across all possible frequency channels ($c_{kk}(t, \tau)$), an can convey
784 tonotopic information through the identity of the recording channel $k$. In order to isolate purely
785 temporal correlation cues, we remove the tonotopic information from the temporal correlation
786 signal by randomly reordering the frequency channels during the classification step.
787

## Sound Database for Auditory Model and Classifier

789     Sounds representing 13 acoustic categories are obtained from a variety of digital sound
790 sources. 195 sounds from 13 different acoustic categories are used to build distributions of

791 dynamic, spectro-temporal correlation statistics of natural/man-made sounds. Sound segments
792 are chosen so that they have minimal background noise and are drawn from 3 broad classes:
793 vocalizations, environmental sounds, and man-made noises. Vocalizations include 1) Single bird
794 songs (various species), 2) Cat meowing (single or multiple cats), 3) Dog barking (single or
795 multiple dogs) and 4) Human speech (male speaker). Environmental sounds include 5) Bird
796 chorus (various species), 6) Speech babble (in various environments e.g. bars, super markets,
797 squares), 7) Fire, 8) Thunder rain, 9) Flowing water (rivers and streams), 10) Wave (ocean/lake
798 waves), 11) Wind. Finally, man-made noises consist of 12) Bell (church or tower bells) and 13)
799 Automobile engines (different vehicles). Each category contains 15 sounds, each 10 s long,
800 sampled at $F_s = 44.1$kHz (see Table S1 for sources and full list of tracks used).
801
802 **Auditory Model**
803    Sounds are analyzed through a cochlear filter bank model of the auditory periphery that
804 decomposes the sound using frequency organized cochlear filters. The cochlear filter banks
805 consists of tonotopically arranged gamma-tone filters (Irino & Patterson, 1996). These filters
806 have a sharp high frequency cutoff and shallow low frequency tails that resemble the tuning
807 functions of auditory nerve fibers. The $k$-th gamma-tone filter has an impulse response function:
808

$$h_k(t) = a_k t^{n-1} e^{-2\pi b_k t} cos(2\pi f_k t + \phi)$$

809
810 where $k$ is the filter channel, $t$ denotes time, $b_k$, and $f_k$ denote the filter bandwidth and center
811 frequency. The filter gain coefficient $a_k$ is chosen so that the filter passband gain is 1; filter order
812 $n$ and filter phase $\phi$ are 3 and 0, respectively. Filter bandwidths are chosen to follow
813 perceptually derived critical bandwidths $b_k = 25 + 75(1 + 1.4 f_k^2)^{0.69}$ [47,48]. We use $L$=58
814 frequency channels with center frequencies $f_k$ ranging from 100 Hz to 16 kHz in 1/8 octave
815 steps. In the first stage of processing, the sound $s(t)$ is passed through the cochlear filterbank
816 model:
817

$$s_k(t) = s(t) * h_k(t)$$

818
819 where $*$ represents the convolution operator. The outputs of the filter bank are next passed
820 through a nonlinear envelope extraction stage, which models the characteristics of the hair cell.
821 We first compute the magnitude of the analytic signal (Cohen,1995):
822

$$s_{A,k}(t) = |s_k(t) + jH\{s_k(t)\}|$$

823
824 where $H\{\cdot\}$ is the Hilbert transform operator and $j = \sqrt{-1}$. Then the temporal envelope for each
825 channel is obtained by convolving the rectified analytic signal magnitude with a b-spline low
826 pass filter with cutoff frequency of 500 Hz (transition bandwidth of 125 Hz and stopband
827 attenuation of 60 dB):
828

$$S_k(t) = S_{A,k}(t) * h_{synapse}(t)$$

829
830 which models the low-pass filtering of the hair cell synapse ($h_{synapse}(t)$). The lowpass filtered
831 envelopes are then down-sampled to 1 kHz for modeling. Here we refer to the time-varying

832    envelopes of the cochlear filters as the cochlear spectrogram and use the notation $S(t, f_k) = S_k(t)$. The cochlear spectrograms, $S(t, f_k)$, of all sound segments are normalized to have zero mean and unit variance.

834    mean and unit variance.

835

## Nonstationary Spectro-temporal Correlation Statistics

837    Since natural sounds are often nonstationary, we measure not just the long-term correlation, but the time-varying or "short-term" correlation statistics between the frequency organized channels in the cochlear model representation. This nonstationary representation is then used to quantify the contribution of the sound correlation statistics to sound categorization. The short-term correlation statistics that we use are similar to those for the neural data analysis except that we use the frequency channels from the cochlear spectrogram rather than the neural signals. The running short-term correlation function $\Phi$ between the cochlear spectrogram channel $k$ and $l$ is computed according to [27]:

845

$$\Phi_{kl}(t, \tau) = \int_{-\infty}^{+\infty} S_k(\gamma) S_l(\gamma - \tau) W^2(t - \gamma) d\gamma$$

846

847    where $W^2(\gamma)$ is a sliding window function that determines the temporal resolution of the correlation measurement, $t$ is the time, $\tau$ is the cross-correlation delay. Here $W^2(\gamma)$ is a Kaiser window ($\beta = 3.4$) where the overall window resolution (corresponding to two standard deviations of the Kaiser window width; varied between 25 to 566 ms in 1/2 octave steps) is varied to quantify the effects of the correlation temporal resolution on categorization performance. The range of permissible cross correlation delays ($-\tau_W$ to $+\tau_W$) correspond to half the window size and is thus varied between -12.5 to +12.5 for the highest resolution to -283 to +283 ms for the coarsest resolution. Conceptually, at each time point, the short-term correlation performs a correlation between the locally windowed envelope signals $S_k(t)W(t - \gamma)$ and $S_l(t - \tau)W(t - \gamma)$ to estimate the localized correlation statistics of the cochlear envelopes.

857    To remove the influence of spectral power on the correlation measurements, the short-term correlations are normalized

859

$$c_{kl}(t, \tau) = \frac{\Phi_{kl}(t, \tau)}{\sqrt{\sigma_k^2(t) \cdot \sigma_l^2(t, \tau)}}$$

860

861    where

862

$$\sigma_k^2(t) = \int_{-\infty}^{+\infty} S_k(\gamma) W(t - \gamma) d\gamma$$

863

$$\sigma_l^2(t, \tau) = \int_{-\infty}^{+\infty} S_l(\gamma - \tau) W(t - \gamma) d\gamma$$

864

865    are the time-varying and delay varying (for channel $l$) power $k$-th and $l$-th spectrogram channels. Again, as with a Pearson correlation coefficient, this short-term spectro-temporal correlation is bounded between -1 and 1.

868  To assess the contribution of temporal and spectral correlations on the sound
869  categorization performance, we perform a secondary analysis where the spectro-temporal
870  correlation function is decomposed into its purely spectral or temporal components, following a
871  similar framework as for the neural data analysis. To evaluate the spectral correlations, we
872  consider only the correlations at $\tau = 0$ (no time lag between different frequency channels). For
873  temporal correlations, $\tau$ ranges from 0 up to half of the Kaiser window length. Only auto-
874  correlation functions or correlation functions of a channel vs. itself are considered in temporal
875  correlation analysis although all possible frequency channels are involved in spectral correlation
876  analysis. Correlations between different frequency channels computed at zero time-lags are
877  referred to as purely spectral correlation components whereas the auto-correlations computed at
878  different time lags for each frequency channel are referred to as temporal correlation throughout.
879

## Stationarity and Ensemble Diversity Indices

881  Given that sounds in our database are quite varied, ranging from isolated vocalizations to
882  environmental sounds consisting of superposition of many individual acoustic events, we seek to
883  characterize the overall degree of stationarity in the short-term correlation statistics for each of
884  the ensemble. Furthermore, since sound recordings are obtained from different sources and
885  animal species (e.g., for vocalizations) all of which could influence the overall category statistics
886  we also seek to quantify the overall diversity of the short-term correlation statistics of each
887  ensemble. When considering sound categorization, we might expect that stationary sounds with
888  minimal diversity across an ensemble would be most easily recognized.
889  For each sound, the sampled short-term spectro-temporal correlation $c_{kl}(t, \tau)$ (computed
890  using $\tau_W$=100 ms) is rearranged and expressed as a time-dependent vector function $\bar{c}(t)$ with
891  dimensions $M \cdot L^2$ at each time point, where $M = 99$ is the number of time lags used for the
892  short-term correlation and $L$ is the number of frequency channels. The stationarity index of each
893  sound is defined and calculated as
894

$$SI = 1 - \frac{\langle \|\bar{c}(t) - \langle \bar{c}(t) \rangle\| \rangle}{\|\langle \bar{c}(t) \rangle\|}$$

895
896  where $\|\cdot\|$ is the vector norm and $\langle \cdot \rangle$ is a time-average. Conceptually, the SI corresponds to the
897  time-average variance of the short-term correlation normalized by the total power of the time-
898  averaged short-term correlation. As such, it measures the average normalized variability across
899  time and is bounded between 0 and 1, where 1 indicates that the moment-to-moment variance of
900  the short-term correlation is 0 thus indicating a high degree of stationarity. By comparison, when
901  the moment-to-moment variance is high the index approaches 0 indicating a highly nonstationary
902  short-term correlation function.
903  We next define and measure the category diversity index (*CDI*), which is designed to
904  measure the degree of homogeneity or heterogeneity in the short-term spectro-temporal
905  correlation functions for each of the 13 sound categories studied. To do so, we first compute the
906  time-average correlation function for each sound in a given ensemble, $\bar{c}_n = \langle \bar{c}_n(t) \rangle$, where
907  $n$=1…15 is an index representing the sounds for each sound category. The *CDI* is defined and
908  computed as
909

$$CDI = \frac{\|\bar{c}_n - E[\bar{c}_n]\|}{\|E[\bar{c}_n]\|}$$

910
911 where $E[\cdot]$ is the expectation operator taken across the sound ensemble (equivalent to an average
912 across sounds, *n)*. Conceptually, the *CDI* corresponds to the variance of the time-average short-
913 term correlation taken across the ensemble of sounds normalized by the power (norm) of the
914 time and ensemble average short-term correlation.
915     For a particular sound category, a *CDI* near zero is indicative of low diversity
916 (homogeneity) such that the short-term spectro-temporal correlation functions of that ensemble
917 are quite similar from sound-to-sound and thus closely resemble the average ensemble
918 correlations. By comparison, an *CDI* of 1 indicates a high degree of heterogeneity (high
919 diversity) so that the short-term spectro-temporal correlations are quite different from sound-to-
920 sound.
921
922 **Dimensionality Reduction and Distribution Model**
923     To reduce the dimensionality of the categorization problem, we use principal component
924 analysis (PCA). For spectral correlation statistics, the entries of the correlation matrix at zero
925 time-lag ($c_{kl}(t,0)$) are considered as features while time points are used as observations or trials.
926 For temporal correlation statistics, the correlations at different time lags within single frequency
927 channels ($c_{kk}(t,\tau)$) are considered features. Both time points and different frequency channels
928 are treated as observations so that temporal information is not specific to any particular
929 frequency channel. For further analysis we use only the highest ranked principal components that
930 explain 90% of the variability in the data (26 PCs for spectral; 8 for temporal; 87 for spectro-
931 temporal).
932     Using the low-dimensional representations of the spectral, temporal, or spectro-temporal
933 correlations, we model the distributions of principal components for each sound category with a
934 Gaussian mixture model (GMM). For each sound category *i* we learn a multivariate probability
935 distribution:
936

$$P(x|m = i) = \sum_{k=1}^{N_c} a_{i,k} N(x; \mu_{i,k}, C_{i,k})$$

937
938 where $x$ is the low-dimensional vector of PCA scores, $m$ are the sound categories, $P(x|m)$ is the
939 multivariate PDF of sound features for category $m$, $a_k \geq 0$ is the weight of the $k$-th mixture
940 component, $N(x; \mu_k, C_k)$ is the PDF of a multivariate normal distributions with mean $\mu_k$ and
941 covariance $C_k$, and $N_c$ is the number of Gaussian components used for modeling. To avoid ill-
942 conditioned covariance matrices, we constrain $C_k$ to be diagonal.
943     In order to find the optimum number of Gaussian components ($N_c$) required to model the
944 data, we compute cross-validated likelihoods for different numbers of Gaussian mixtures ($N_c$=1
945 to 20). The optimal $N_c$ values are 5, 8 and 13 for temporal, spectral, and spectro-temporal,
946 respectively.
947
948 **Bayesian Classifier**
949     Given the mixture model for each sound, we then use a Bayesian classifier for sound
950 identification. The features fed to the classifier consist of the principal component scores from
951 either the sound's short-term spectral, temporal or spectro-temporal correlation. Since we are
952 interested in how categorization performance change with the sound duration, we consider

953 feature vectors $x = [x_1, x_2, ..., x_N]$ consisting of principal component scores obtained from
954 successive, windowed segments of sounds at different time samples ($t_{1...N}$) selected so that
955 adjacent sound segments do not overlap. We then evaluate the posterior probability of each
956 sound segment under the different mixture models. The most probable case is chosen according
957 to the maximum a posteriori (MAP) decision rule:
958

$$\hat{m}_{MAP}(x) = \underset{m}{argmax}\, P(m|x)$$

959

960 In practice, we find the MAP category by using Bayes rule and maximizing the log-likelihood.
961 We assume that the categories are equi-probable a priori and that the features at each time
962 sample are conditionally independent. Thus, the MAP category is obtained by maximizing:
963

$$log(P(m|x)) = log\left[P(m)\prod_{n=1}^{N}\frac{P(x_n|m)}{P(x_n)}\right] \propto \sum_{n=1}^{N} log(P(x_n|m)).$$

964
965

**Cross Validation**

966

967    To avoid over-fitting, we use a leave-one-out cross-validation, where all sounds are used
968 to build the model distributions except for one sound which is used for testing (Larose & Larose,
969 2015). Because there is only one sound validated, each validation iteration produces a 0 % or 100
970 % correct classification rate. The procedure is repeated iteratively over all sounds and the
971 average performance is obtained as the average classification rate across all iterations. The total
972 sound duration used for the validation is varied by selecting $N$ consecutive time window
973 segments as described above from each sound under the test to be categorized; the selected $N$-
974 segments start from the very beginning of the sound up to the end of the sound. Values of $N$ are
975 varied in 1/2 octave steps starting with $N$=1 up to the maximum value allowed by the sound
976 duration.

977

**Optimal Temporal Resolution and Integration Time for Categorizing Sounds**

978

979    As previously demonstrated for auditory neurons an optimal temporal resolution can be
980 identified for neural discrimination of natural sounds, and neural discrimination performance
981 improves with the increasing sound duration [41,42]. For this reason, we seek to identify both the
982 optimal temporal resolution that maximizes categorization performance as well as the integration
983 time of the sound classifier. The temporal resolution of the correlation signals is varied by
984 changing the sliding window temporal resolution, $\tau_W$, between 25-566 ms in 1/2 octave steps.
985 Classification performance curves vary with $\tau_W$, exhibiting concave behavior with a clear
986 maximum that is used to identify the optimal window time-constant. The classifier performance
987 also increased in an approximately exponential fashion with the overall sound duration. The
988 classifier performance also increases with the overall sound duration. The classifier integration
989 rise-time, $\tau_c$, is defined as the amount of time required to achieve 90% of the asymptotic
990 performance measured at 10 sec duration.

991
992

993

994
995     REFERENCES
996
997     1       Escabi, M. A. & Schreiner, C. E. Nonlinear spectrotemporal sound analysis by neurons in
998             the auditory midbrain. *J Neurosci* **22**, 4114-4131, doi:20026325
999     22/10/4114 [pii] (2002).
1000    2       McDermott, J. H. & Simoncelli, E. P. Sound texture perception via statistics of the
1001            auditory periphery: evidence from sound synthesis. *Neuron* **71**, 926-940, doi:S0896-
1002            6273(11)00562-9 [pii]
1003    10.1016/j.neuron.2011.06.032 (2011).
1004    3       Rodriguez, F. A., Chen, C., Read, H. L. & Escabi, M. A. Neural modulation tuning
1005            characteristics scale to efficiently encode natural sound statistics. *J Neurosci* **30**, 15969-
1006            15980, doi:10.1523/JNEUROSCI.0966-10.2010 (2010).
1007    4       Singh, N. C. & Theunissen, F. E. Modulation spectra of natural sounds and ethological
1008            theories of auditory processing. *J Acoust Soc Am* **114**, 3394-3411 (2003).
1009    5       Kulkarni, A. & Colburn, H. S. Role of spectral detail in sound-source localization. *Nature*
1010            **396**, 747-749 (1998).
1011    6       Oxenham, A. J., Bernstein, J. G. & Penagos, H. Correct tonotopic representation is
1012            necessary for complex pitch perception. *Proceedings of the National Academy of*
1013            *Sciences of the United States of America* **101**, 1421-1425, doi:10.1073/pnas.0306958101
1014            (2004).
1015    7       Geffen, M. N., Gervain, J., Werker, J. F. & Magnasco, M. O. Auditory perception of self-
1016            similarity in water sounds. *Front Integr Neurosci* **5**, 15, doi:10.3389/fnint.2011.00015
1017            (2011).
1018    8       Gervain, J., Werker, J. F. & Geffen, M. N. Category-specific processing of scale-invariant
1019            sounds in infancy. *PLoS ONE* **9**, e96278, doi:10.1371/journal.pone.0096278 (2014).
1020    9       Elhilali, M., Ma, L., Micheyl, C., Oxenham, A. J. & Shamma, S. A. Temporal coherence in
1021            the perceptual organization and cortical representation of auditory scenes. *Neuron* **61**,
1022            317-329, doi:S0896-6273(08)01053-2 [pii]
1023    10.1016/j.neuron.2008.12.005 (2009).
1024    10      Barlow, H. in *Sensory Communication*      (MIT Press, 1961).
1025    11      Chechik, G. *et al.* Reduction of information redundancy in the ascending auditory
1026            pathway. *Neuron* **51**, 359-368, doi:S0896-6273(06)00512-5 [pii]
1027    10.1016/j.neuron.2006.06.030 (2006).
1028    12      Nelken, I., Rotman, Y. & Bar Yosef, O. Responses of auditory-cortex neurons to structural
1029            features of natural sounds. *Nature* **397**, 154-157 (1999).
1030    13      Hsu, A., Woolley, S. M., Fremouw, T. E. & Theunissen, F. E. Modulation power and phase
1031            spectrum of natural sounds enhance neural encoding performed by single auditory
1032            neurons. *J Neurosci* **24**, 9201-9211, doi:10.1523/JNEUROSCI.2449-04.2004 (2004).
1033    14      Natan, R. G., Carruthers, I. M., Mwilambwe-Tshilobo, L. & Geffen, M. N. Gain Control in
1034            the Auditory Cortex Evoked by Changing Temporal Correlation of Sounds. *Cereb Cortex*
1035            **27**, 2385-2402, doi:10.1093/cercor/bhw083 (2017).
1036    15      Miller, L. M. & Schreiner, C. E. Stimulus-based state control in the thalamocortical
1037            system. *J Neurosci* **20**, 7011-7016 (2000).

1038 16 deCharms, R. C. Information coding in the cortex by independent or coordinated
1039 populations. *Proc Natl Acad Sci U S A* **95**, 15166-15168 (1998).
1040 17 Downer, J. D., Niwa, M. & Sutter, M. L. Task engagement selectively modulates neural
1041 correlations in primary auditory cortex. *J Neurosci* **35**, 7565-7574,
1042 doi:10.1523/JNEUROSCI.4094-14.2015 (2015).
1043 18 Chen, C., Read, H. L. & Escabi, M. A. Precise feature based time-scales and frequency
1044 decorrelation lead to a sparse auditory code
1045 . *J Neurosci* **32**, 8454-8468 (2012).
1046 19 Cohen, M. R. & Newsome, W. T. Estimates of the contribution of single neurons to
1047 perception depend on timescale and noise correlation. *J Neurosci* **29**, 6635-6648,
1048 doi:10.1523/JNEUROSCI.5179-08.2009 (2009).
1049 20 Abbott, L. F. & Dayan, P. The effect of correlated variability on the accuracy of a
1050 population code. *Neural Comput* **11**, 91-101 (1999).
1051 21 Jeffress, L. A. A place theory of sound localization. *J Comp Physiol Psychol* **41**, 35-39
1052 (1948).
1053 22 Shamma, S. & Klein, D. The case of the missing pitch templates: how harmonic
1054 templates emerge in the early auditory system. *J Acoust Soc Am* **107**, 2631-2644 (2000).
1055 23 Langner, G. Neural processing and representation of periodicity pitch. *Acta Otolaryngol*
1056 *Suppl* **532**, 68-76 (1997).
1057 24 Fu, Q. J. & Shannon, R. V. Recognition of spectrally degraded and frequency-shifted
1058 vowels in acoustic and electric hearing. *J Acoust Soc Am* **105**, 1889-1900 (1999).
1059 25 Greenberg, S. Speaking in shorthand – A syllable-centric perspective for understanding
1060 pronunciation variation. *Speech Communication* **29**, 159-176 (1999).
1061 26 Khatami, F., Read, H. L., Wöhr, M. & Escabi, M. A. Origins of scale invariance in
1062 vocalization sequences and speech 4. *PLoS computational biology* (in press).
1063 27 Sayers, B. M. & Cherry, E. C. Mechanisms of Binaural Fusion in the Hearing of Speech. *J.*
1064 *Acoust. Soc. Am.* **29**, 973-987 (1957).
1065 28 Shannon, R. V., Zeng, F. G., Kamath, V., Wygonski, J. & Ekelid, M. Speech recognition
1066 with primarily temporal cues. *Science* **270**, 303-304 (1995).
1067 29 Smith, Z. M., Delgutte, B. & Oxenham, A. J. Chimaeric sounds reveal dichotomies in
1068 auditory perception. *Nature* **416**, 87-90, doi:10.1038/416087a (2002).
1069 30 Churchland, A. K. *et al.* Variance as a signature of neural computations during decision
1070 making. *Neuron* **69**, 818-831, doi:10.1016/j.neuron.2010.12.037 (2011).
1071 31 Salinas, E. & Sejnowski, T. J. Correlated neuronal activity and the flow of neural
1072 information. *Nat Rev Neurosci* **2**, 539-550, doi:10.1038/35086012 (2001).
1073 32 Cohen, M. R. & Maunsell, J. H. Attention improves performance primarily by reducing
1074 interneuronal correlations. *Nat Neurosci* **12**, 1594-1600, doi:10.1038/nn.2439 (2009).
1075 33 Oliver, D. L. & Morest, D. K. The central nucleus of the inferior colliculus in the cat. *J*
1076 *Comp Neurol* **222**, 237-264 (1984).
1077 34 Read, H. L., Winer, J. A. & Schreiner, C. E. Modular organization of intrinsic connections
1078 associated with spectral tuning in cat auditory cortex. *Proc Natl Acad Sci U S A* **98**, 8042-
1079 8047, doi:10.1073/pnas.131591898
1080 98/14/8042 [pii] (2001).

1081    35    Storace, D. A., Higgins, N. C. & Read, H. L. Thalamic label patterns suggest primary and
1082          ventral auditory fields are distinct core regions. *J Comp Neurol* **518**, 1630-1646,
1083          doi:10.1002/cne.22345 (2010).
1084    36    Rabinowitz, N. C., Willmore, B. D., King, A. J. & Schnupp, J. W. Constructing noise-
1085          invariant representations of sound in the auditory pathway. *PLoS biology* **11**, e1001710,
1086          doi:10.1371/journal.pbio.1001710 (2013).
1087    37    Fischer, B. J., Pena, J. L. & Konishi, M. Emergence of multiplicative auditory responses in
1088          the midbrain of the barn owl. *J Neurophysiol* **98**, 1181-1193 (2007).
1089    38    Joris, P. X., Smith, P. H. & Yin, T. C. Coincidence detection in the auditory system: 50
1090          years after Jeffress. *Neuron* **21**, 1235-1238 (1998).
1091    39    McDermott, J. H., Schemitsch, M. & Simoncelli, E. P. Summary statistics in auditory
1092          perception. *Nat Neurosci* **16**, 493-498, doi:nn.3347 [pii]
1093    10.1038/nn.3347 (2013).
1094    40    Shadlen, M. N. & Newsome, W. T. Neural basis of a perceptual decision in the parietal
1095          cortex (area LIP) of the rhesus monkey. *J Neurophysiol* **86**, 1916-1936 (2001).
1096    41    Narayan, R., Grana, G. & Sen, K. Distinct time scales in cortical discrimination of natural
1097          sounds in songbirds. *J Neurophysiol* **96**, 252-258, doi:01257.2005 [pii]
1098    10.1152/jn.01257.2005 (2006).
1099    42    Engineer, C. T. *et al.* Cortical activity patterns predict speech discrimination ability. *Nat*
1100          *Neurosci* **11**, 603-608, doi:nn.2109 [pii]
1101    10.1038/nn.2109 (2008).
1102    43    Schnupp, J. W., Garcia-Lazaro, J. A. & Lesica, N. A. Periodotopy in the gerbil inferior
1103          colliculus: local clustering rather than a gradient map. *Front Neural Circuits* **9**, 37,
1104          doi:10.3389/fncir.2015.00037 (2015).
1105    44    Roark, R. M. & Escabí, M. A. B-spline design of maximally flat and prolate spheroidal-
1106          type FIR filters. *IEEE Transactions on Signal Processing* **47**, 701-716 (1999).
1107    45    Rodriguez, F. A., Read, H. L. & Escabi, M. A. Spectral and temporal modulation tradeoff
1108          in the inferior colliculus. *J Neurophysiol* **103**, 887-903, doi:10.1152/jn.00813.2009
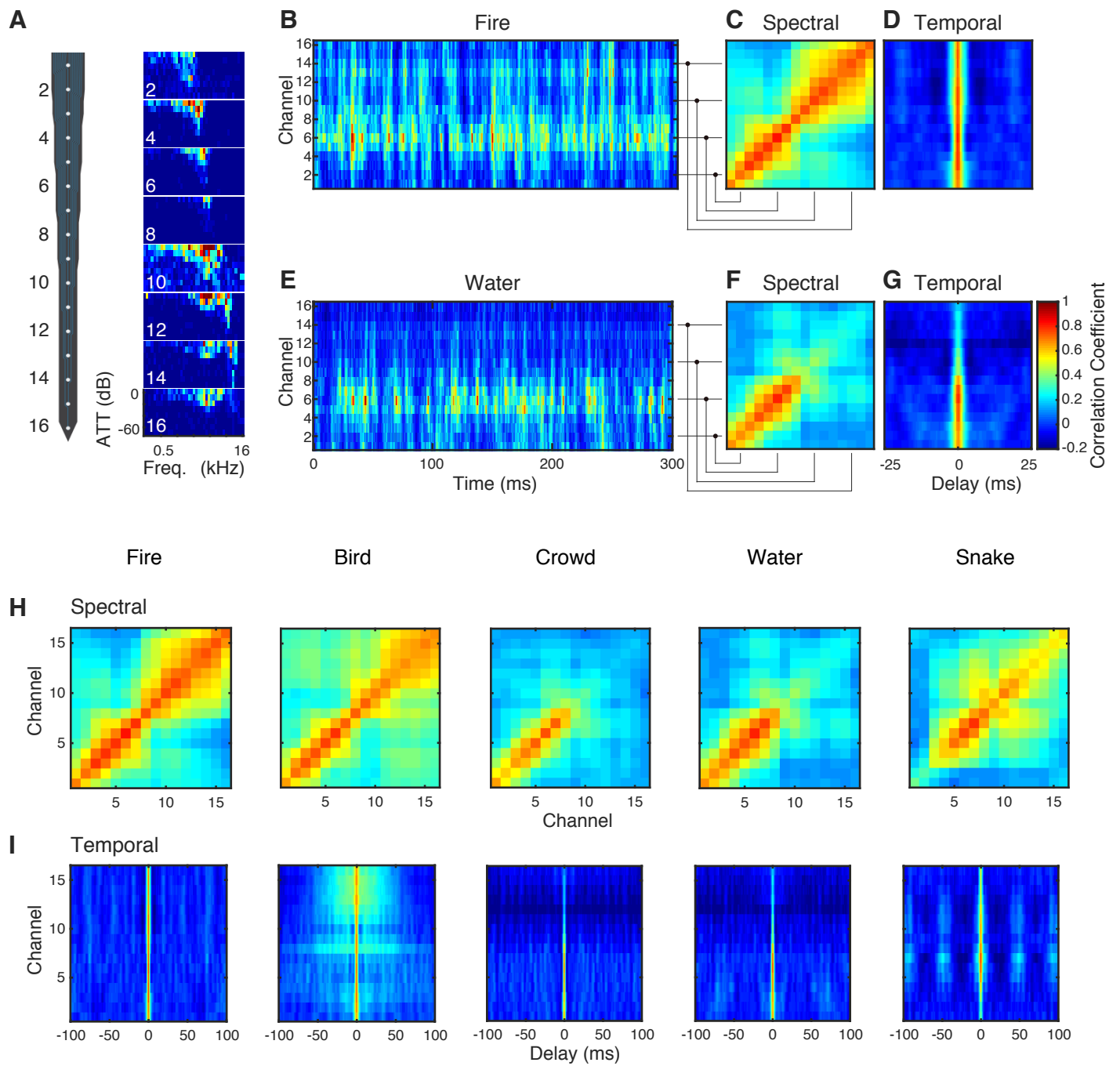1109          (2010).
1110    46    Zheng, Y. & Escabi, M. A. Distinct roles for onset and sustained activity in the neuronal
1111          code for temporal periodicity and acoustic envelope shape. *J Neurosci* **28**, 14230-14244,
1112          doi:28/52/14230 [pii]
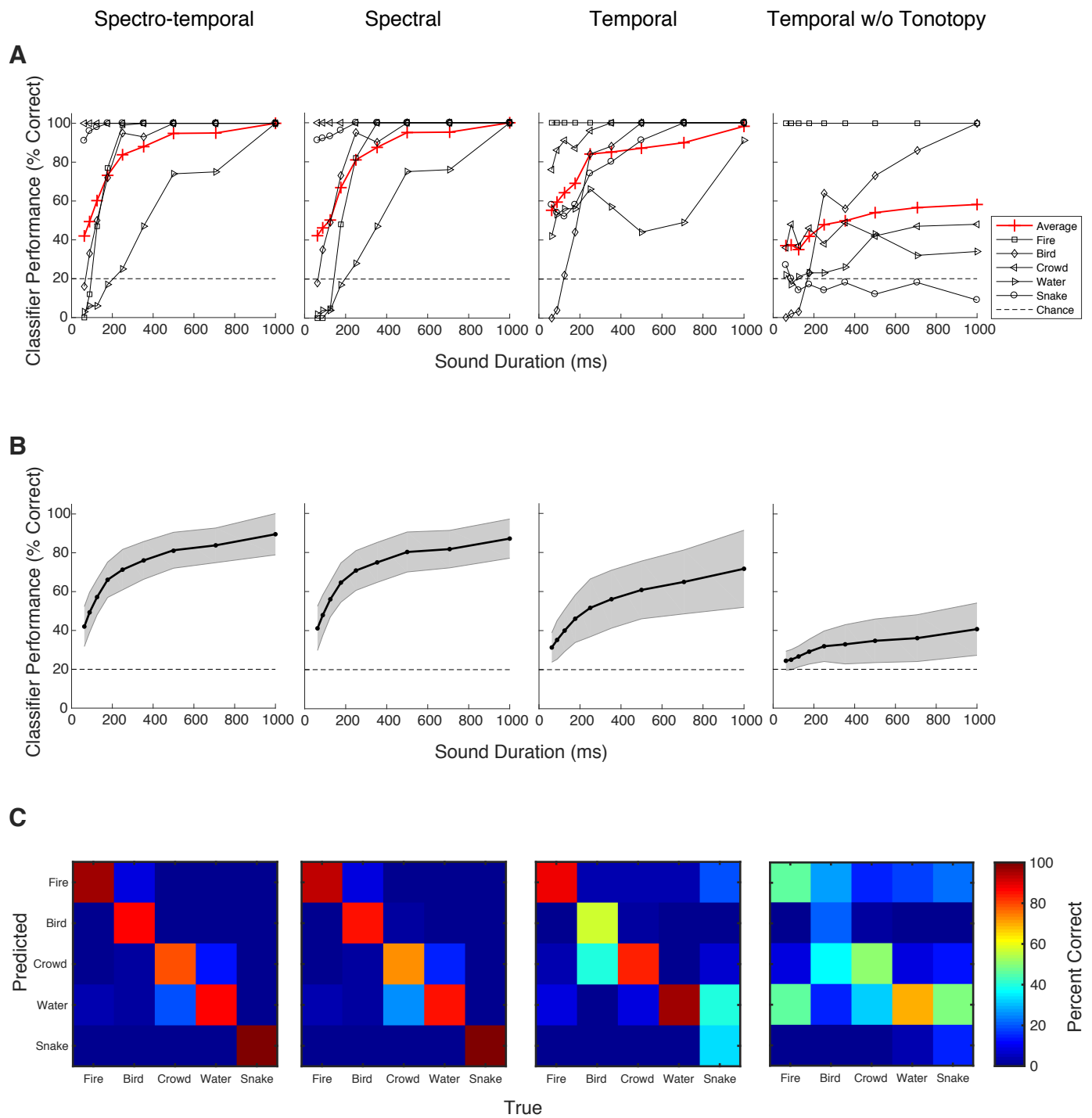1113    10.1523/JNEUROSCI.2882-08.2008 (2008).
1114    47    Fletcher, H. Auditory patterns. *Rev Mod Phys* **12**, 47-65 (1940).
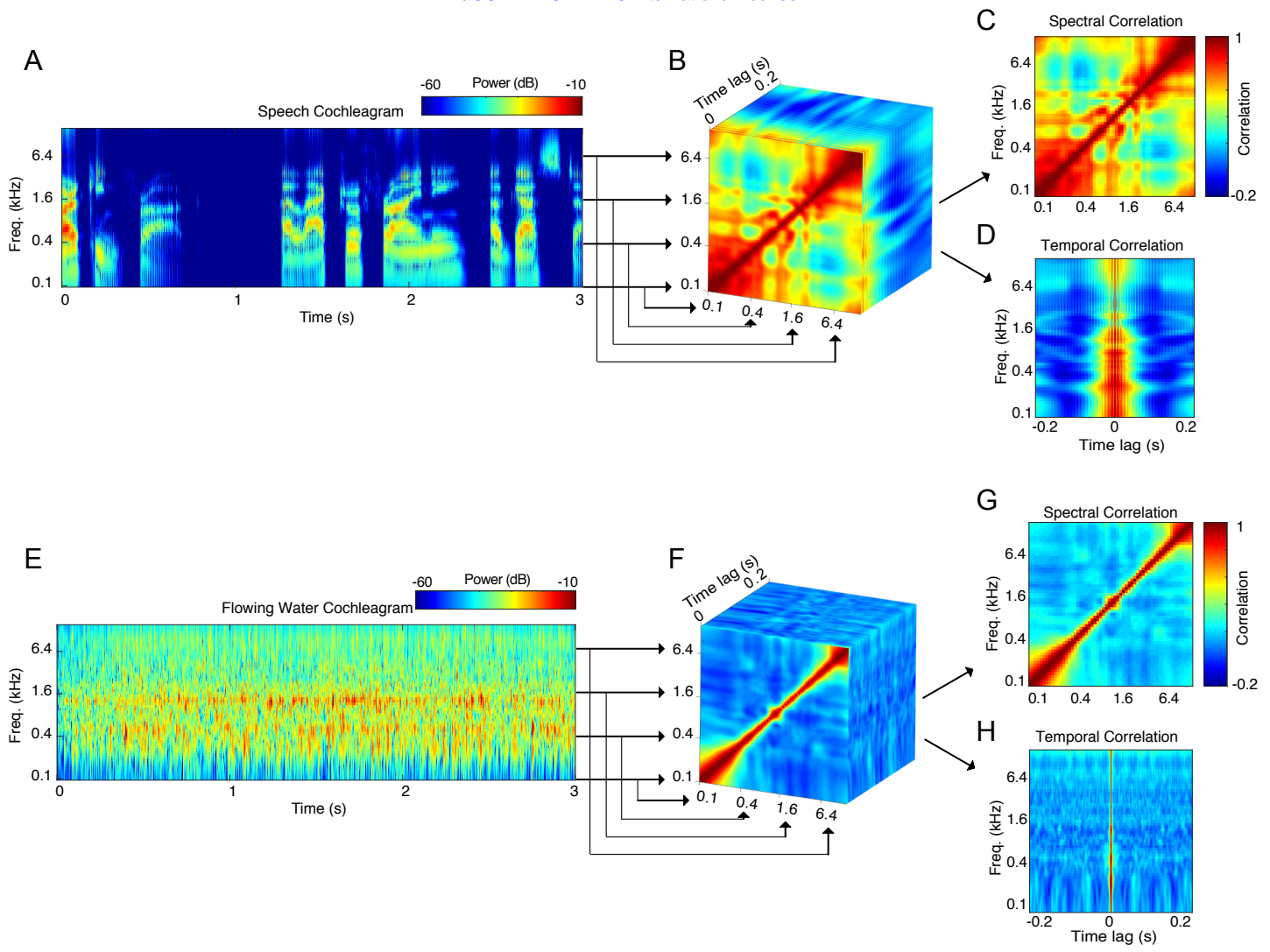1115    48    Zwicker, E., Flottorp, G. & Stevens, S. S. Critical band width in loudness summation. *J*
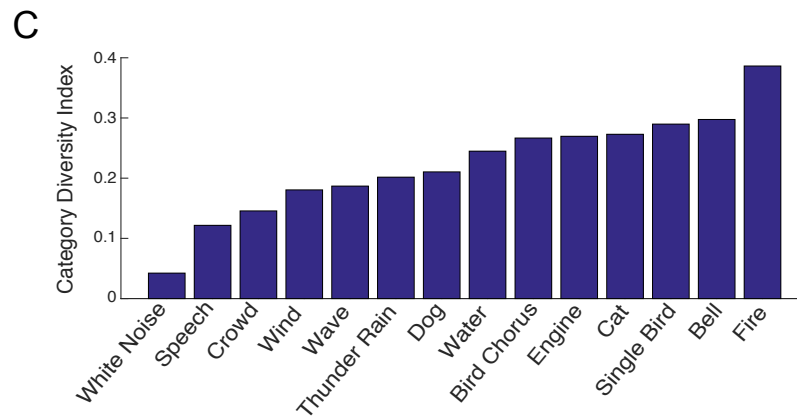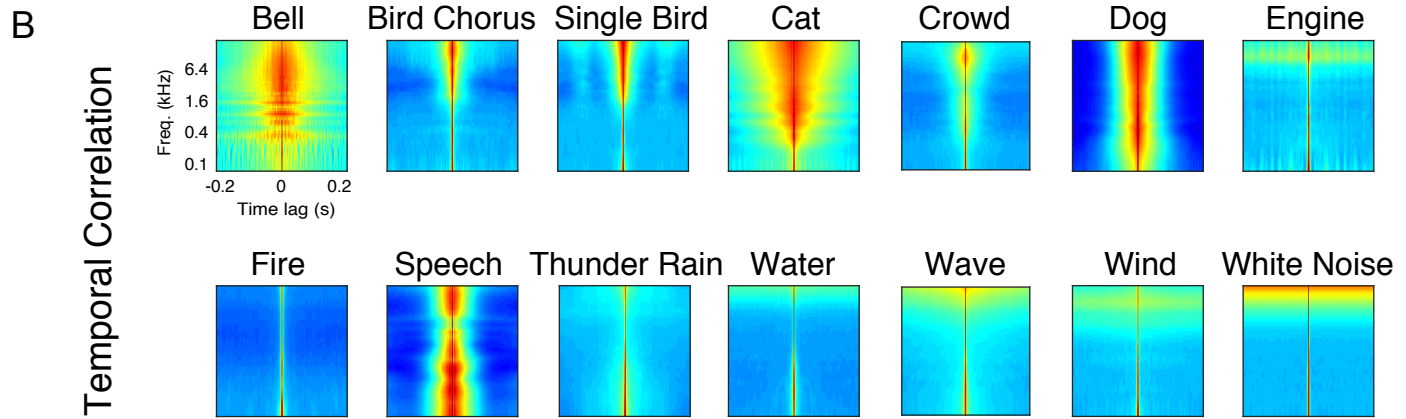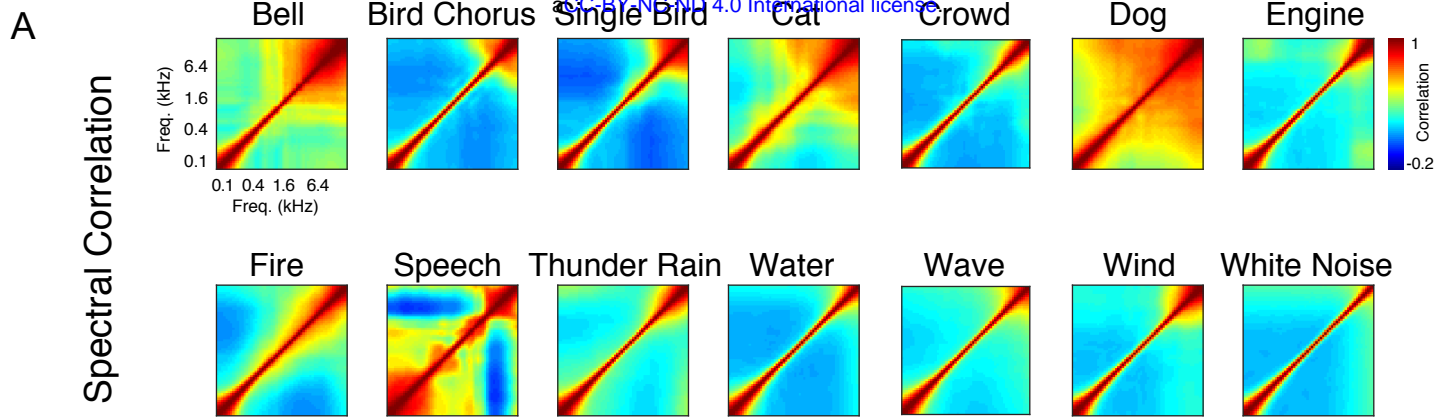1116          *Acoust Soc Am* **29**, 548-557 (1957).
1117    49    Cohen, M. R. & Kohn, A. Measuring and interpreting neuronal correlations. *Nat Neurosci*
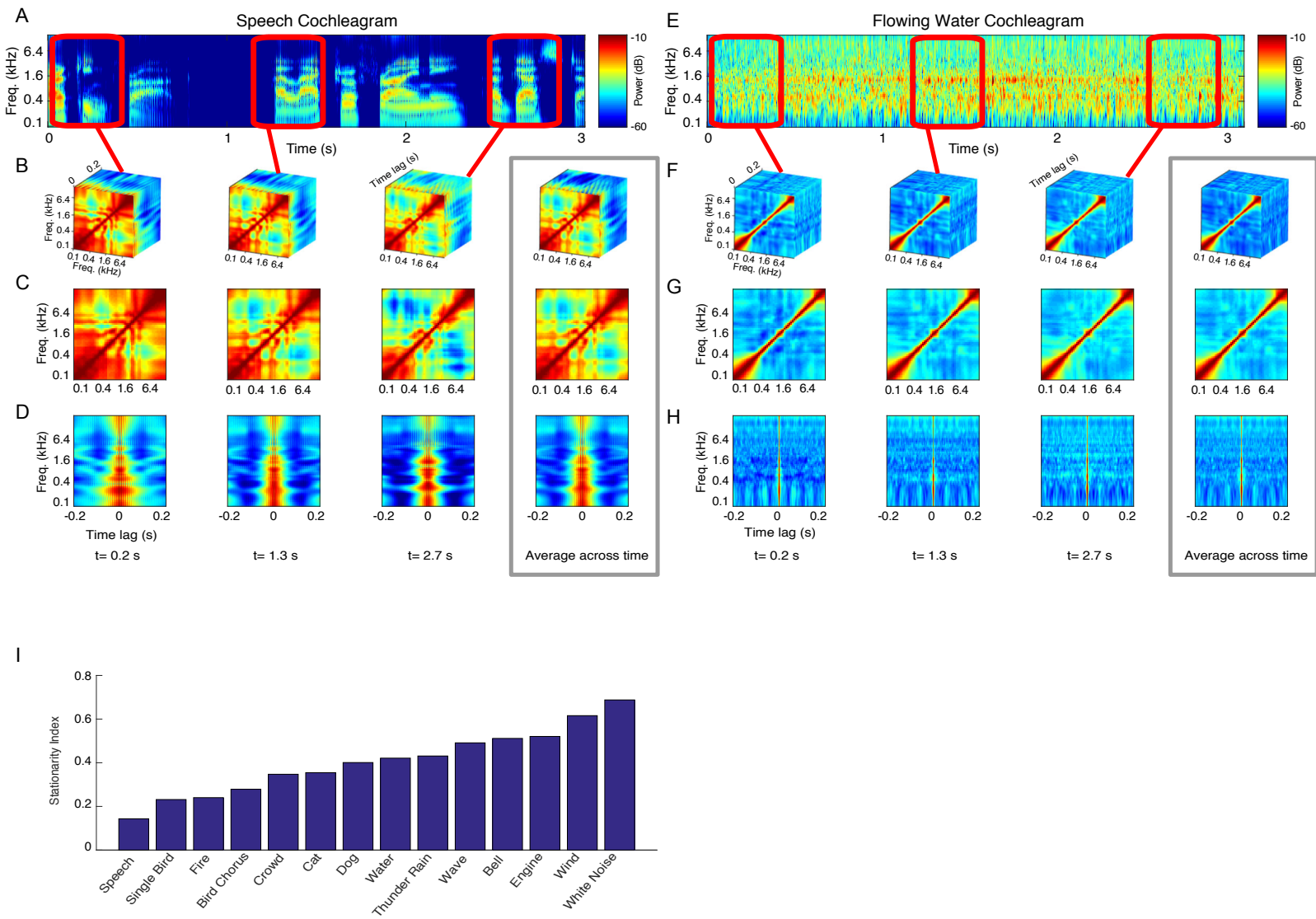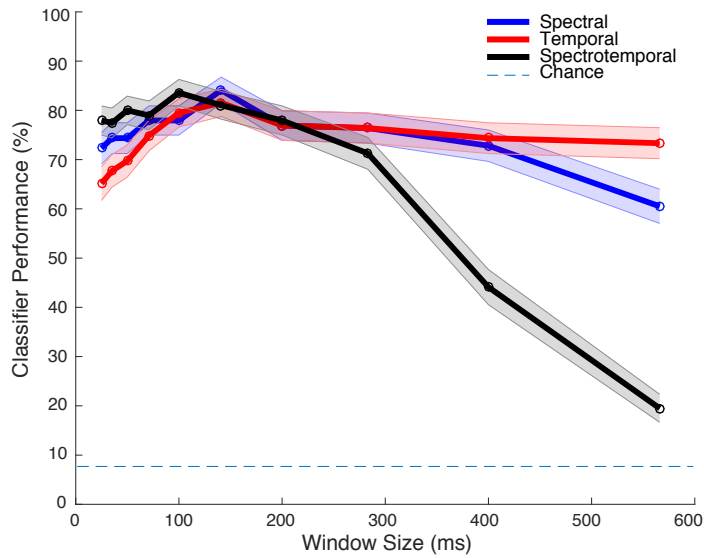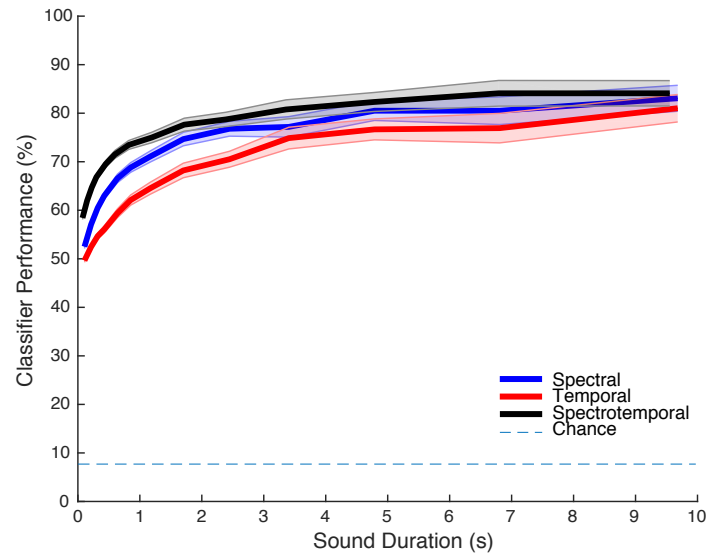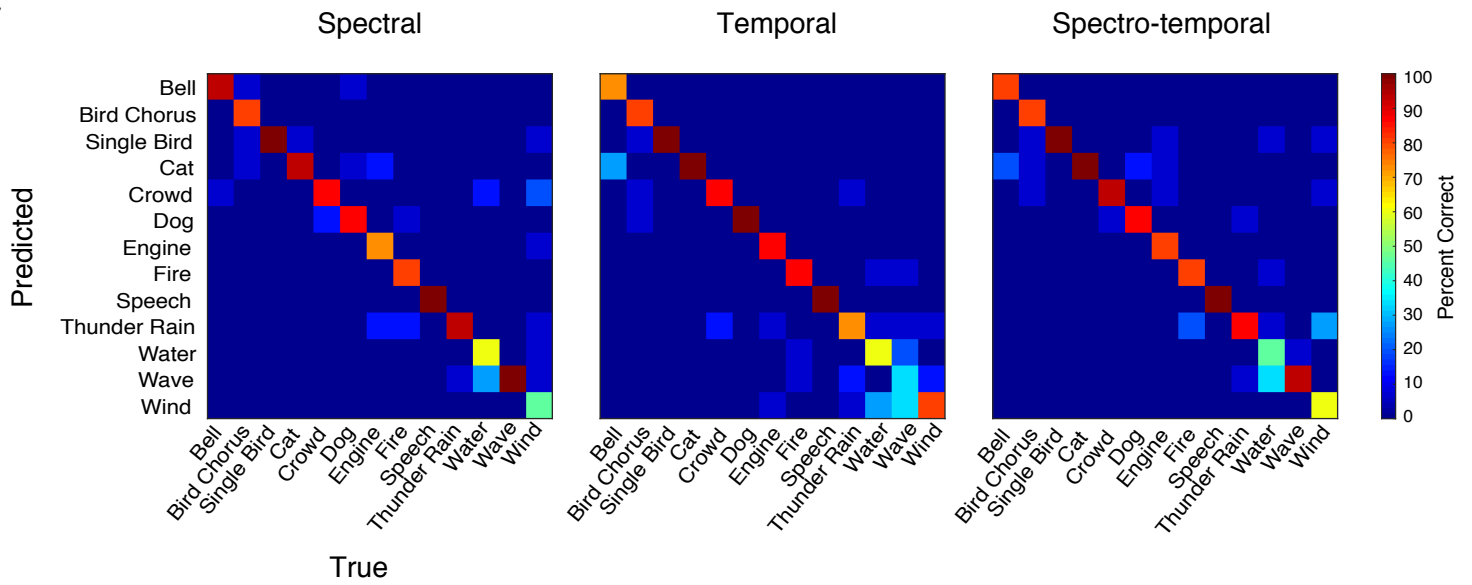1118          **14**, 811-819, doi:10.1038/nn.2842 (2011).
1119