

Recovering sound sources from embedded repetition

Josh H. McDermott^{a,1}, David Wroblewski^b, and Andrew J. Oxenham^b

^aCenter for Neural Science, New York University, New York, NY 10003 and ^bDepartment of Psychology, University of Minnesota, Minneapolis, MN 55455

Edited by Edward Adelson, Massachusetts Institute of Technology, Cambridge, MA, and approved December 1, 2010 (received for review April 8, 2010)

Cocktail parties and other natural auditory environments present organisms with mixtures of sounds. Segregating individual sound sources is thought to require prior knowledge of source properties, yet these presumably cannot be learned unless the sources are segregated first. Here we show that the auditory system can bootstrap its way around this problem by identifying sound sources as repeating patterns embedded in the acoustic input. Due to the presence of competing sounds, source repetition is not explicit in the input to the ear, but it produces temporal regularities that listeners detect and use for segregation. We used a simple generative model to synthesize novel sounds with naturalistic properties. We found that such sounds could be segregated and identified if they occurred more than once across different mixtures, even when the same sounds were impossible to segregate in single mixtures. Sensitivity to the repetition of sound sources can permit their recovery in the absence of other segregation cues or prior knowledge of sounds, and could help solve the cocktail party problem.

auditory scene analysis | cocktail party problem | generative models of sound | natural sound statistics | sound segregation

Auditory scenes generally contain multiple sources, the sounds from which add together to produce a mixed signal that enters the ears. In most behavioral contexts, however, it is the sources, not the mixture, that are of interest. This is often termed the “cocktail party problem”—organisms must infer individual sound sources from ambiguous mixtures of sounds (1–7).

Recovering individual sound sources from an auditory scene requires assumptions, or priors, about what sources are like (8). For instance, listeners implicitly assume that frequency components that are regularly spaced (9, 10), begin and end simultaneously (11), or have similar distributions of binaural spatial cues (12) belong to the same sound. Listeners also use knowledge of specific familiar sound classes, filling in masked syllable segments in ways that are consistent with known speech acoustics (13).

Priors on sounds are thus used by the auditory system and must somehow be acquired; yet natural environments rarely feature isolated sound sources from which they could be readily learned. Organisms face a “chicken and egg” problem—sound sources must be separated from mixtures for their properties to be learned, but to separate sources from mixtures, listeners need to know something about their characteristics to begin with.

It is possible that priors are at least partially built into the auditory system by evolution, or that listeners can learn them from occasionally hearing sound sources in isolation. In this paper we consider an alternate, complementary, solution—that listeners might detect sources as repeating spectro-temporal patterns embedded in the acoustic input. Both individual sound sources and their mixtures produce combinations of acoustic features, but because mixtures result from multiple independent sources, the feature configurations that they produce are unlikely to occur repeatedly with consistency. Repetition is thus a signature of individual sources. The repetition of a sound source is generally not explicit in the signal that enters the ear, due to the corruption of a source’s acoustic signature by other sounds. However, repeating sources induce temporal regularities in the mixed auditory input, which we suggest are detected and used by the auditory system to recover sound sources.

To explore this idea, we studied the conditions under which listeners could identify novel sound sources that they only ever heard in mixtures with other sounds. We developed a method to

synthesize novel sounds that shared some of the correlation structure of natural sounds (14–16) but that lacked strong grouping cues, and presented listeners with mixtures of these sounds. Listeners were generally unable to identify the sounds composing a single such mixture, but when presented with multiple mixtures of a particular target sound with various others, they heard the target repeating across mixtures and could reliably identify it. Even two presentations of the target yielded a significant benefit.

Our results indicate that listeners detect latent repeating spectro-temporal structure within sound mixtures and from this can identify individual sound sources. Sound source repetition thus serves as a powerful cue that can “bootstrap” performance in situations in which other bottom-up cues and top-down knowledge are unavailable, and as such may play an important role in auditory scene analysis.

Results

Generative Model for Sounds. To test whether source repetition might by itself be sufficient for sound segregation, it was important both to use novel sounds, so that familiarity would not enable segregation, and to minimize the presence of bottom-up grouping cues in our test stimuli. However, we wanted our results to have real-world relevance, and thus to use stimuli with some similarity to natural sounds. We met these goals by modeling the time-frequency decomposition (spectrogram) of a sound as a Gaussian-distributed random variable with correlations that resembled those in natural sounds.

We first generated spectrograms for sets of spoken words (Fig. 1A) and animal vocalizations (Fig. 1B). Such spectrograms generally share a simple property: the energy at nearby points tends to be similar (14–16). This is evident when the correlation between pairs of spectrogram cells is plotted as a function of their time and frequency offset (Fig. 1C and D). For both classes of natural sounds, correlations are high for small offsets and decline with separation in time or frequency, whereas for noise signals they are absent. Such results follow from the common finding that natural modulation spectra (related to correlation functions via the Fourier transform) peak at low modulation frequencies (14–16) and thus exhibit correlations over moderate time/frequency scales.

We used correlation functions similar to those of natural sound sets (Fig. 1C and D) to generate a covariance matrix, each element of which was the covariance between two spectrogram cells. Spectrograms were drawn from the resulting Gaussian distribution and applied to samples of white noise, yielding novel sounds (Fig. 1E and F). Related stimuli result from constraining the modulation spectrum of noise (16); our spectrogram-domain method had advantages in implementing our task (*SI Materials and Methods*). Although our stimuli shared important statistical properties of real sounds, they lacked the grouping cues provided by abrupt temporal onsets and harmonic spectral structure, both of which are important for sound segregation (1, 2) but which are not captured by second-order correlations.

Author contributions: J.H.M. and A.J.O. designed research; J.H.M. and D.W. performed research; J.H.M. analyzed data; and J.H.M. and A.J.O. wrote the paper.

The authors declare no conflict of interest.

This article is a PNAS Direct Submission.

¹To whom correspondence should be addressed. E-mail: jhm@cns.nyu.edu.

This article contains supporting information online at www.pnas.org/lookup/suppl/doi:10.1073/pnas.1004765108/-DCSupplemental.

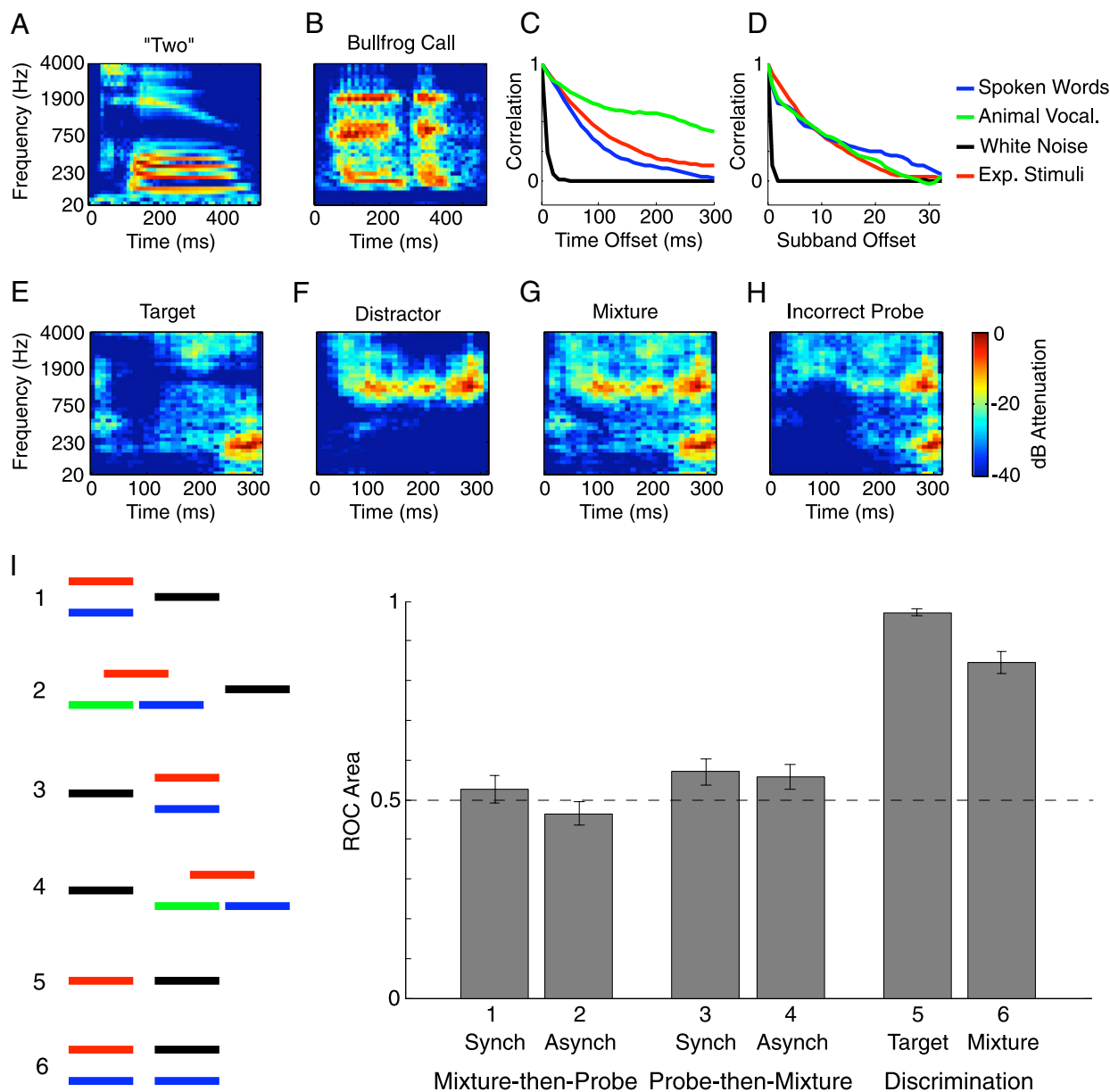


Fig. 1. Stimulus generation and results of Experiment 1. (A and B) Time-frequency decomposition of a spoken word and a bullfrog vocalization. (C and D) Correlation between nearby time-frequency cells as a function of their temporal (C) and spectral (D) separation. (E and F) Two spectrograms generated by our model. (G) Spectrogram of the mixture of the sounds from E and F. (H) Spectrogram of an incorrect probe sound, generated to be physically consistent with the mixture in G. (I) Results and stimulus configurations from Experiment 1. Line segments represent sounds; sounds presented simultaneously are drawn as vertically displaced. Distinct sounds are indicated by different colors. Red segments represent target sounds, and black segments represent probe sounds. Error bars denote SEs. The dashed line represents the chance performance level.

Performance-Based Measure of Sound Segregation. We assessed sound segregation by presenting mixtures of sounds (Fig. 1G) followed by a probe sound. Listeners judged whether the probe had been present in the mixture(s). The probe was either one of the sounds in the mixture(s), termed the “target” sound, or another sound with statistics similar to the target (Fig. 1H). In the latter case, the probe was constrained to be physically consistent with the mixture (such that—like the target—it never had more energy than the mixture). Each target was presented only once per experiment, so that subjects could not learn the targets from the probes.

Following the probe presentation, subjects selected one of four responses (“sure no,” “no,” “yes,” or “sure yes”) to indicate whether they thought the probe was one of the sounds in the mixture. These responses were used to generate a receiver operating characteristic (ROC) curve. The area beneath the curve

was our performance measure (17); chance and perfect performance corresponded to areas of 0.5 and 1, respectively. All of the effects reported here are evident in the stimulus examples available at http://www.cns.nyu.edu/~jhm/source_repetition.

Experiment 1: Sound Segregation with Single Mixtures. We began by presenting subjects with single mixtures of two sounds (Fig. 1I). Sound segregation should permit a listener to judge whether a subsequent probe sound was one of the sounds in the mixture. However, performance was generally at chance levels, even after considerable practice [condition 1: $t(9) = 0.64, P = 0.54$]. Performance remained close to chance when we included a third sound and made the sounds asynchronous [condition 2: $t(9) = -0.65, P = 0.53$]. Asynchrony should enhance the bottom-up grouping cue provided by onset differences between sources (1, 2, 11);

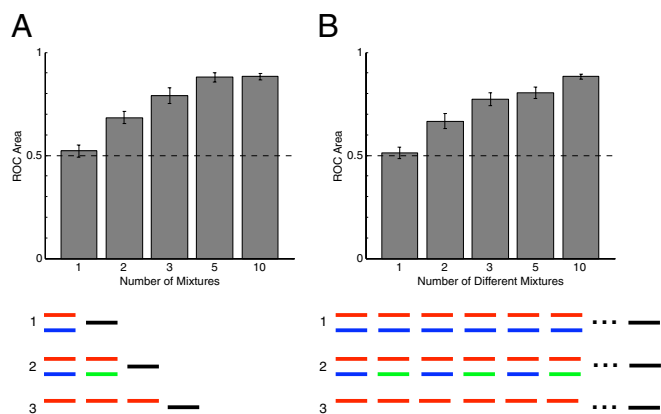


Fig. 2. Effect of multiple mixtures on sound source recovery. (A) Different numbers of mixtures were presented. (B) Ten mixtures were presented in all conditions, and the number of different mixtures was varied. Conventions here and elsewhere are as in Fig. 1*i*. Red segments represent target probes, black segments represent incorrect probes, and different colors represent different sounds. Schematics for conditions with 5 and 10 mixtures are omitted.

the lack of effect suggests that any onsets in our stimuli were too weak to support segregation. We also tried presenting the probe sound before the mixture, so that subjects knew what sound to listen for, but performance was still not significantly different from chance [condition 3, synchronous: $t(9) = 2.23, P = 0.053$; condition 4, asynchronous: $t(9) = 1.8, P = 0.1$], although there was a small effect of hearing the probe first [$F(1,9) = 7.33, P = 0.02$].

The poor performance was not due to an inability to discriminate different synthetic sounds; the correct and incorrect probe sounds were easily distinguished when presented in isolation [condition 5: $t(9) = 56.1, P < 10^{-12}$]. Moreover, when the target and incorrect probe sounds for a particular mixture were each mixed with the same unrelated second sound, the resulting mixtures themselves were discriminable [condition 6: $t(9) = 12.5, P < 10^{-6}$]. Thus, chance performance in the sound segregation task was not due to limits on encoding of the mixtures (as it would be if the stimulus differences needed to perform our task were completely masked by the other sound in the mixture). Rather, performance was evidently limited by the inability to segregate the mixture into two sounds. The subjective experience of listening to the mixtures was consistent with this conclusion. The mixtures usually sounded like a single sound that was qualitatively different from the target sound.

These results indicate that our stimuli met our principal objectives. Despite having some naturalistic structure, they lacked the grouping cues needed to segregate them from a mixture. This made them well suited to our primary goal of testing whether sound structure could be extracted from multiple occurrences of a target sound.

Experiment 2: Sound Segregation with Multiple Mixtures. To test whether listeners could benefit from sound source repetition across mixtures, we presented target sounds repeatedly, each time mixed with a different “distractor” sound. Despite the difficulty of segregating single mixtures, a target presented more than once in succession was usually heard repeating through the mixtures, and listeners rapidly developed an impression of it. In Experiment 2a we quantified this benefit, varying the number of mixtures and measuring how well subjects could discriminate correct from incorrect target probes. Performance was again at chance levels with a single mixture, but improved as subjects heard more mixtures (Fig. 2*A*). Performance was significantly improved even with two mixtures [$t(9) = 3.66, P = 0.005$] and appeared to asymptote with about five mixtures.

To rule out the possibility that the improvement with multiple mixtures was due merely to repeated exposure to the target, in Experiment 2b we held the number of mixtures constant at 10,

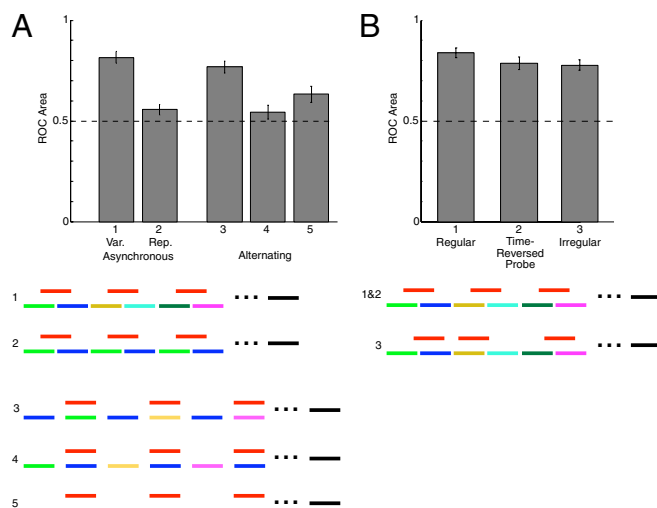


Fig. 3. Stimuli and results of Experiment 3. (A) Effect of mixture variability persists with asynchronous and alternating presentation. Conditions 3 and 4 differ in the pairing of the target with variable (condition 3) or repeated (condition 4) distractors. (B) Subjects can perform task even when incorrect probes are time-reversed versions of the target sound, or when the target sound is presented irregularly.

but varied how many different mixtures occurred in the sequence. In the single-mixture condition, subjects heard the same mixture 10 times. The 10-mixture condition was the same as in Experiment 2a. The other conditions repeatedly presented two, three, or five mixtures in a fixed order over the course of the sequence, with each mixture containing the target sound.

Performance again steadily increased with the number of different mixtures (Fig. 2*B*), even though the target was always presented the same number of times. The ability to hear the target sound thus appears to depend on the number of different mixtures that a listener hears, not on the total number of target presentations. An ANOVA comparing the two experiments showed a main effect of the number of different mixtures [$F(4,36) = 115.35, P < 0.0001$], but no effect of experiment type [$F(1,9) = 0.73; P = 0.42$] and no interaction [$F(4,36) = 0.59, P = 0.67$]. See *SI Results* for additional controls.

As with the single mixtures of Experiment 1, the sounds composing the single repeated mixtures tended to blend together and rarely bore close resemblance to the target sound. This is consistent with the idea that listeners detect repeating sound structure and attribute it to individual sources; when the same mixture repeats, it is heard as a source, and the target structure is no more apparent than when it is heard only once.

Experiment 3a: Asynchronous Mixtures. Experiment 2 featured synchronously presented sounds, but distinct sources in real-world scenes are generally asynchronous. Experiment 3a confirmed that the benefit of multiple distinct mixtures persisted when the target and distractors were temporally offset to better resemble natural conditions (Fig. 3*A, Left*, condition 1 vs. condition 2). As before, a single repeated mixture yielded near-chance performance, but presenting different mixtures in succession enabled discrimination of the target sound [$F(1,7) = 116.87, P < 0.0001$]. The effect of multiple mixtures in this case swamps that of any grouping cue provided by the asynchrony (consistent with the weak onsets in our sounds), and is not specific to synchronously presented sounds.

The effect was also evident when the target sound was presented with every other distractor in a sequence (Fig. 3*A, Right*, conditions 3–5). When the distractors that co-occurred with the target varied (condition 3), performance was well above chance, even though the distractors that alternated with the target repeated ($P = 0.004$, sign test). But when the distractor sequence

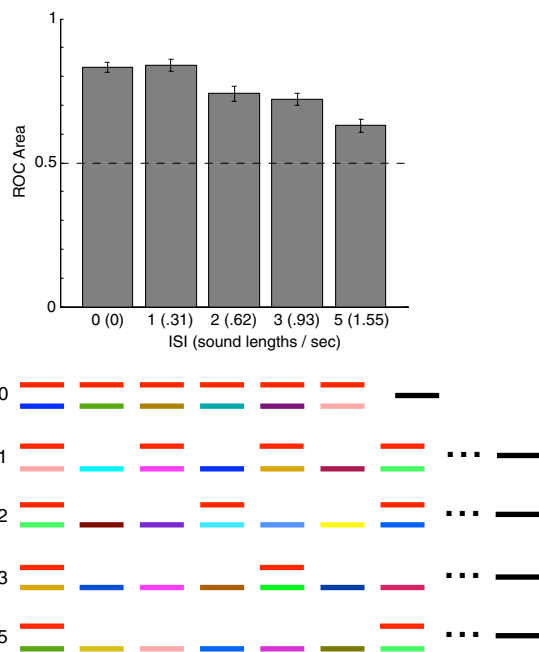


Fig. 4. Effect of interstimulus interval. In all conditions, the target sounds (shown in red) were presented six times. Condition 0 is identical to the variable mixture conditions of Experiment 2 except for the number of target presentations.

was phase-shifted by a target length, so that the repeating distractors co-occurred with the target (condition 4), the target was generally unidentifiable. When every distractor repeated (condition 5), performance tended to be intermediate between the other two conditions (significantly worse than the variable condition and better than the repeated condition, $P = 0.008$ and 0.06 , respectively, sign test; also better than condition 1, Experiment 2b, $P = 0.06$). This configuration is reminiscent of some used in studies of pure tone streaming (18). In this condition, the repetition of the distractor may compete with that of the mixture.

Experiment 3b: Spectrotemporal Structure and Irregular Presentation. To test whether listeners extracted the temporal structure of sounds in addition to their spectral content, in Experiment 3b we presented variable mixtures but used a time-reversed version of the target sound for the incorrect probe (that thus had the same power spectrum as the target but differed in temporal structure). As shown in Fig. 3B, performance remained high when distinguishing between the correct and the time-reversed probes, although there was a slight advantage with our standard incorrect probes [$F(2,18) = 4.03$, $P = 0.04$]. Listeners thus derived a spectrotemporal profile for the target sound and did not merely encode the average spectrum of the mixture sequence. Performance also remained high when the targets were presented at irregular temporal intervals (Fig. 3B), indicating that periodically occurring acoustic structure was not necessary for the effect ($P < 0.0001$ for both conditions, two-tailed t test).

Experiment 4: Temporal Integration. If the benefit of multiple mixtures on sound segregation reflects the extraction of repeating structure from the auditory input, it should be constrained by the short-term storage capacity of the auditory system; to recognize that a structure repeats, the input must be stored over the repetition time. We examined the effect of target spacing on subjects' ability to extract the target from a mixture sequence, holding the number of target presentations fixed at six but varying how frequently the targets occurred (Fig. 4). Performance was unaffected by short delays but declined steadily thereafter [$F(4,24) = 22.98$, $P < 0.0001$]. The results are consistent with an integration process

that tracks acoustic structure using an auditory memory buffer, although they leave open the question of whether time delays or the intervening acoustic input are driving the effect. Either way, it appears that when the storage capacity of the integration process is exceeded, repetition becomes difficult to track.

Computational Schemes for Extracting Embedded Repetition. It is easy to envision simple computational schemes in which the structure of a repeating source could be extracted from mixtures. As a proof of concept, Fig. 5 illustrates one such approach. A target estimate is initialized to the first segment of the mixture sequence and over time is refined through an averaging process that is time-locked to peaks in the cross-correlation of the target estimate and the spectrogram (*SI Methods*). The correlation peaks reveal the delay at which the signal contains the target, and the averaging (taking the pointwise minimum of the previous target estimate and the current spectrogram segment) combines information across mixtures. Although the estimation process is constrained by the averaging window (*SI Methods*), it does not require knowledge of the target duration, repetition pattern, or other characteristics.

Fig. 5 shows a spectrogram of a sequence of mixtures of a target sound with various others (A), followed by spectrograms of a sequence of target estimates derived for this mixture sequence (B), graphs showing the cross-correlation between each successive target estimate and the next 700-ms block of the spectrogram (C), and a spectrogram of the true target (D). The correlation peaks occur at the onset of the target in the mixture, and the estimation process converges on the true target after several iterations (see also *SI Results*, Experiment 6).

Discussion

The recovery of individual sound sources from mixtures of multiple sounds is a central challenge of hearing. Our results suggest one solution: a sound source can be recovered if it occurs more than once and is not always mixed with the same other sounds. This is true even in cases where other grouping cues are impoverished to the point that a single instance of the source is unsegmentable. The auditory system evidently detects repeating spectro-temporal structure embedded in mixtures, and interprets this structure as a sound source. Repetition of sound sources is not explicit in the input to the ear, because the source waveform is generally corrupted at each presentation by other sounds. Source repetition can nonetheless be detected by integrating information over time. Listeners in our experiments were able to form detailed impressions of sound sources that they only ever heard in mixtures, and thus were able to recover this latent structure.

Source repetition can be viewed as another acoustic grouping cue, but it is distinct from other cues in one important respect—its use does not require prior knowledge of sound characteristics. Other grouping cues are rooted in particular properties of natural sounds (e.g., the “bottom-up” cues of common onset or harmonicity) or attributes specific to individual sounds or sound classes (e.g., the “top-down” cues of speech acoustics). Such properties serve as cues because they characterize the particular sorts of sounds found in the world. Knowledge of these sound properties thus must first be internalized by the auditory system from the environment, either over the course of evolution or by learning during an organism's development. Repetition, in contrast, requires only the assumption that sound sources maintain some consistency over time. Our finding that repetition alone can support segregation suggests that it can bootstrap the auditory system in situations where characteristics of sound sources are not yet known, be it early in development or in unfamiliar auditory environments.

The practical utility of this phenomenon for sound segregation obviously depends on the presence of repeating sounds. Not all sounds occur repetitively, but repetition is nonetheless common to natural auditory environments. Examples include the sounds of rhythmic motor behaviors (e.g., walking, running, scratching, clapping) and repetitive physical processes (e.g., branches swaying, water trickling). It is also striking that many animal

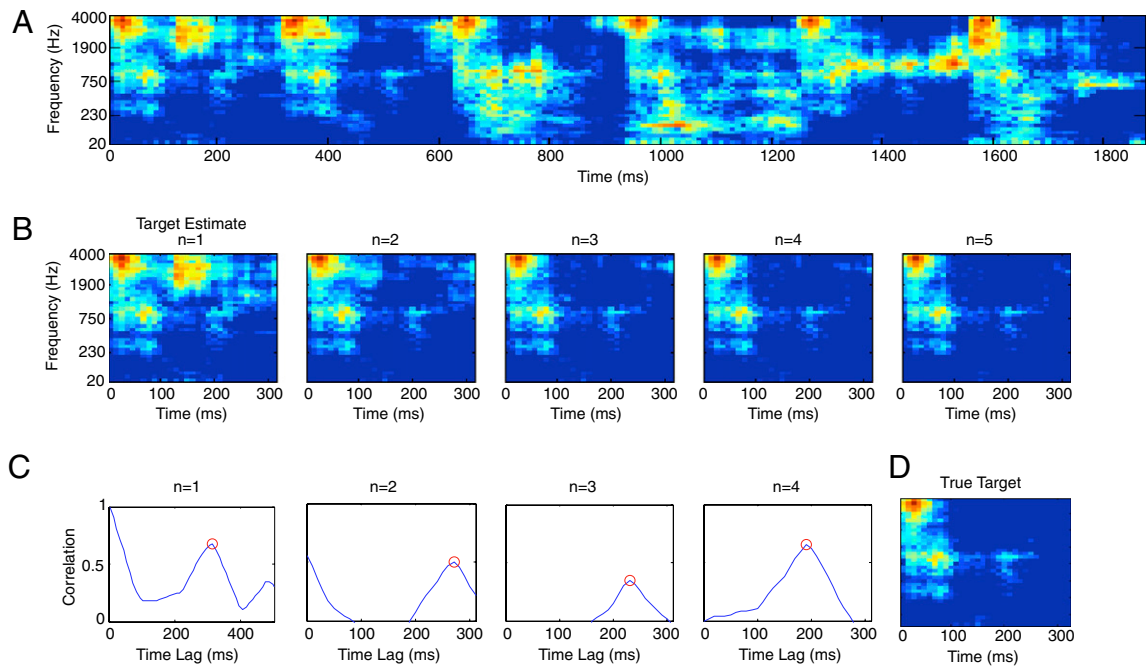


Fig. 5. A candidate computational scheme to extract a repeating target sound from mixtures. (A) Spectrogram of a sequence of mixtures of one target sound with various distractors. (B) Spectrograms of target sound estimates after each iteration of the algorithm. Only the first 300 ms is shown for ease of comparison with D. (C) Cross-correlation of target estimate with the next block of the input spectrogram from A, as a function of the time shift applied to the spectrogram block. The red circle denotes the peak of the correlation function as found by a peak-picking algorithm. (D) Spectrogram of the true target sound. Note the resemblance to the target estimate after five iterations, shown directly above.

vocalizations consist of repetitions of a single call (19), and as such would benefit from repetition-based segregation. Although the targets in our experiments repeated exactly, we found informally that moderate variation in the exemplars had little effect on the ability to hear the target repeating. This is not surprising from a computational standpoint; if the repeating sounds produce a peak in the correlation function, as they will when their variation is not excessive, then an algorithm like that of Fig. 5 will recover their central tendency. “Fingerprinting” techniques for detecting repeating patterns (20) are an alternative model for repetition detection, and these are particularly tolerant of variability. It thus seems likely that source repetition could play an important role in everyday hearing.

The effect of repetition can be viewed as an extension of Bregman’s “old-plus-new” idea (1), whereby frequencies added to a spectrum are segregated from those that are continuously present. Our effects involve continuity only at an abstract level, because our stimuli had dynamic spectra and were often separated by short gaps (Figs. 3 and 4). Our results thus implicate a mechanism that can extract dynamic spectrotemporal structure (e.g., as in Fig. 5) distinct from the spectral subtraction mechanisms often posited (1). The upshot of this is that repetition can drive the segregation of complex, quasi-realistic sounds from mixtures.

The effects described in this paper are examples of “streaming” (1, 18, 21, 22), in that the repeating targets segregate from the distractors over time. Perhaps because we presented temporally overlapping sounds, our effects differ in some respects from the well-known case of alternating tones that segregate when repeated. We found that sounds segregated only when one of the sounds varied, not when both were repeated. Our findings bear a closer resemblance to the classic finding that repeating tones are easier to detect when accompanying masker tones vary from one presentation to the next (23–25). Those effects are conceptually similar to ours, but the acoustics are considerably different, as are the conditions under which the effects hold. For instance, the tone effects depend on spectral separation between the target tone and the masker, perhaps relying on spectral separation as a bottom-up

segregation cue, and are adversely affected by even brief gaps between tones (25). These differences from our phenomena raise the possibility of distinct mechanisms; the tone effects seem closely related to Bregman’s old-plus-new phenomena, and could have a similar explanation. There is also some conceptual similarity between our results and demonstrations that infants and adults can learn repeating patterns in streams of phonemes (26). This latter case seems likely to represent a distinct phenomenon, given that the patterns are acquired over longer time scales and usually are not consciously accessible.

Our study highlights the experimental use of generative models of sound. Studies of the cocktail party problem have traditionally used unnatural synthetic stimuli (9, 27, 28) or familiar real-world sounds such as speech (3, 10, 12, 29). Generative models have the advantage of producing novel stimuli that lack the confounding effects of familiarity but that share properties of natural sounds. The statistics captured by our model are but a small subset of those characterizing the full distribution of natural sounds, but they nonetheless have two important consequences. First, stimuli with naturalistic modulation are sparse in the time-frequency domain, and thus they do not uniformly mask one another (8). Detection of repetition likely requires some degree of sparsity in the sensory input, because otherwise there would be little to gain from hearing sounds in multiple mixtures; most sounds would mask one another over most of their extent. Second, natural statistics allowed the generation of many stimuli that did not all sound the same. Presumably because the auditory system is tuned to the properties of natural sounds (30–33), in this case spectro-temporal modulation (34), naturalistic stimuli are better discriminated than unnatural stimuli (35). Different samples of white noise, for instance, sound much less distinct than do different samples from our model, which likely would make the task of discriminating targets prohibitively difficult.

Consistent with these notions, pilot experiments with alternative correlation functions indicated that the phenomena do not depend sensitively on their exact shape, but that large deviations from natural correlations do render the stimuli less discriminable and less sparse, to the point that the task becomes impossible. For

instance, we found that the task could not be performed when the stimuli were different samples of white noise. Although repetition of individual samples of white noise is sometimes noticeable (36, 37), their perceptual similarity and spectrotemporal uniformity apparently precludes this when samples are embedded in mixtures. It thus was important to use a naturalistic sound model. Sparsity is likely crucial to the phenomenon, and the discriminability of natural stimuli facilitated the experimental task.

The utility of source repetition could extend to vision and olfaction, which also confront scene analysis problems. Organisms receive multiple overlapping objects or odors as sensory input, and repetition might enable the recovery of individual objects or odors without prior knowledge of their characteristics. The problems are not analogous in all of their details (e.g., odors are not defined by their temporal structure, and visual objects do not combine linearly when forming an image, due to occlusion; ref. 7), but the same general principle may apply: a particular mixture of sources (objects or odors) is unlikely to occur repeatedly, such that repeating patterns in the input are diagnostic of single sources. Repeating patterns should induce input correlations that could guide temporal integration and reveal single objects or odors, just as we found with sound.

The cocktail party problem has been believed to be solved via the combination of grouping cues derived from statistical regularities of natural sounds, and knowledge of specific sounds or sound classes. Using a simple generative model to produce novel sounds, we found that sound source repetition provides a third source of information with which to parse sound mixtures, one that the auditory system can use even when other segregation cues are unavailable, and which could perhaps be used to learn other grouping cues. The auditory system seems attuned to repetition, and can use it to succeed in conditions that would otherwise be insurmountable.

Materials and Methods

Sound analysis and synthesis used spectrograms specifying the logarithm of the rms amplitude in a set of time-frequency windows. Spectrograms were generated by first passing a signal through an auditory filter bank, then passing each filter output through a set of time windows. The rms level of the windowed signal yielded the value of a spectrogram cell. Adjacent filters and time windows overlapped by 50%.

Correlations between pairs of spectrogram cells were measured for the initial 500-ms segment of each natural sound. These correlations were averaged across pairs of cells with the same time or frequency offset to yield temporal and spectral correlation functions for each stimulus set, as displayed in Fig. 1 C and D.

Synthetic stimuli with similar correlations were created by modeling the spectrogram as a multivariate Gaussian variable, specified by a mean spectrogram, and a covariance matrix containing the covariance between every pair of spectrogram cells. The mean of each spectrogram cell was set proportional to the corresponding filter bandwidth. The covariance matrix was generated from exponentially decaying correlation functions that approximated the shape of correlation functions for natural sounds. For each pair of cells, the covariance was the product of the corresponding temporal and spectral correlations and a constant variance.

To generate sounds, a time-frequency decomposition was generated for a sample of white noise. The signal in each window was scaled to set its log-amplitude to that of the corresponding cell in a spectrogram sampled from our generating distribution. The results were passed through the filter bank again (as in other analysis and synthesis decompositions; ref. 38) and summed to generate a sound signal. Because adjacent filters and time windows overlapped and thus interfered with each other when amplitudes were altered, the spectrogram of the resulting sound generally differed from the sampled spectrogram from which the sound was generated. However, these differences were subtle, and the intended correlation structure remained present in the sounds, as can be seen in the correlations measured in the synthetic sounds (Fig. 1 C and D).

Methods are described in more detail in *SI Methods*.

ACKNOWLEDGMENTS. We thank Bart Anderson, Stephen David, Monty Escabi, Heather Read, and Nathan Witthoft for comments on the manuscript, and Andrew Schwartz and Barbara Shinn-Cunningham for useful discussions. This work was supported by National Institutes of Health Grant R01 DC 07657.

- Bregman AS (1990) *Auditory Scene Analysis: The Perceptual Organization of Sound* (MIT Press, Cambridge, MA).
- Darwin CJ, Carlyon RP (1995) Auditory grouping. *The Handbook of Perception and Cognition*, ed Moore BCJ (Academic, New York), Vol 6.
- Bronkhorst AW (2000) The cocktail party phenomenon: A review of research on speech intelligibility in multiple-talker conditions. *Acustica* 86:117–128.
- Narayan R, et al. (2007) Cortical interference effects in the cocktail party problem. *Nat Neurosci* 10:1601–1607.
- Bee MA, Micheyl C (2008) The cocktail party problem: What is it? How can it be solved? And why should animal behaviorists study it? *J Comp Psychol* 122:235–251.
- Ehlihalil M, Shamma SA (2008) A cocktail party with a cortical twist: How cortical mechanisms contribute to sound segregation. *J Acoust Soc Am* 124:3751–3771.
- McDermott JH (2009) The cocktail party problem. *Curr Biol* 19:R1024–R1027.
- Ellis DPW (2006) Model-based scene analysis. *Computational Auditory Scene Analysis: Principles, Algorithms, and Applications*, eds Wang D, Brown GJ (Wiley, Hoboken, NJ), pp 115–146.
- Roberts B, Brunstrom JM (1998) Perceptual segregation and pitch shifts of mistuned components in harmonic complexes and in regular inharmonic complexes. *J Acoust Soc Am* 104:2326–2338.
- de Cheveigne A, Kawahara H, Tsuzaki M, Aikawa K (1997) Concurrent vowel identification. I: Effects of relative amplitude and F0 difference. *J Acoust Soc Am* 101:2839–2847.
- Darwin CJ, Ciocca V (1992) Grouping in pitch perception: Effects of onset asynchrony and ear of presentation of a mistuned component. *J Acoust Soc Am* 91:3381–3390.
- Best V, Ozmeral E, Gallun FJ, Sen K, Shinn-Cunningham BG (2005) Spatial unmasking of birdsong in human listeners: Energetic and informational factors. *J Acoust Soc Am* 118:3766–3773.
- Warren RM (1970) Perceptual restoration of missing speech sounds. *Science* 167:392–393.
- Voss RF, Clarke J (1975) “1/f noise” in music and speech. *Nature* 258:317–318.
- Attias H, Schreiner CE (1997) Temporal low-order statistics of natural sounds. *Advances in Neural Information Processing*, eds Mozer M, Jordan M, Petsche T (MIT Press, Cambridge, MA), Vol 9.
- Singh NC, Theunissen FE (2003) Modulation spectra of natural sounds and ethological theories of auditory processing. *J Acoust Soc Am* 114:3394–3411.
- MacMillan NA, Creelman CD (1991) *Detection Theory: A User's Guide* (Cambridge Univ Press, New York).
- Bregman AS, Pinker S (1978) Auditory streaming and the building of timbre. *Can J Psychol* 32:19–31.
- Wiley RH, Richards DG (1978) Physical constraints on acoustic communication in the atmosphere: Implications for the evolution of animal vocalizations. *Behav Ecol Sociobiol* 3:69–94.
- Cotton C, Ellis D (2009) Finding similar acoustic events using matching pursuit and locality-sensitive hashing. *Proceedings of the IEEE Workshop on Applications of Signal Processing to Audio and Acoustics, New Paltz, New York* (Institute of Electrical and Electronics Engineers, New York).
- Moore BCJ, Gockel H (2002) Factors influencing sequential stream segregation. *Acta Acustica* 88:320–332.
- Snyder JS, Alain C (2007) Toward a neurophysiological theory of auditory stream segregation. *Psychol Bull* 133:780–799.
- Kidd G, Jr., Mason CR, Deliwala PS, Woods WS, Colburn HS (1994) Reducing informational masking by sound segregation. *J Acoust Soc Am* 95:3475–3480.
- Kidd G, Jr., Mason CR, Richards VM (2003) Multiple bursts, multiple looks, and stream coherence in the release from informational masking. *J Acoust Soc Am* 114:2835–2845.
- Micheyl C, Shamma SA, Oxenham AJ (2007) Hearing out repeating elements in randomly varying multitone sequences: a case of streaming?. *Hearing: From Sensory Processing to Perception*, eds Kollmeier B, et al. (Springer, Berlin).
- Saffran JR, Aslin RN, Newport EL (1996) Statistical learning by 8-month-old infants. *Science* 274:1926–1928.
- Carlyon RP (1991) Discriminating between coherent and incoherent frequency modulation of complex tones. *J Acoust Soc Am* 89:329–340.
- McDermott JH, Oxenham AJ (2008) Spectral completion of partially masked sounds. *Proc Natl Acad Sci USA* 105:5939–5944.
- Culling JF, Darwin CJ (1993) Perceptual separation of simultaneous vowels: Within and across-formant grouping by F0. *J Acoust Soc Am* 93:3454–3467.
- Smith EC, Lewicki MS (2006) Efficient auditory coding. *Nature* 439:978–982.
- Nelken I, Rotman Y, Bar Yosef O (1999) Responses of auditory-cortex neurons to structural features of natural sounds. *Nature* 397:154–157.
- Garcia-Lazaro JA, Ahmed B, Schnupp JW (2006) Tuning to natural stimulus dynamics in primary auditory cortex. *Curr Biol* 16:264–271.
- Escabi MA, Miller LM, Read HL, Schreiner CE (2003) Naturalistic auditory contrast improves spectrotemporal coding in the cat inferior colliculus. *J Neurosci* 23:11489–11504.
- Chi T, Gao Y, Guyton MC, Ru P, Shamma SA (1999) Spectro-temporal modulation transfer functions and speech intelligibility. *J Acoust Soc Am* 106:2719–2732.
- Woolley SM, Fremouw TE, Hsu A, Theunissen FE (2005) Tuning for spectro-temporal modulations as a mechanism for auditory discrimination of natural sounds. *Nat Neurosci* 8:1371–1379.
- Kaernbach C (2004) The memory of noise. *Exp Psychol* 51:240–248.
- Agus TR, Thorpe SJ, Pressnitzer D (2010) Rapid formation of robust auditory memories: Insights from noise. *Neuron* 66:610–618.
- Crochiere RE, Webber SA, Flanagan JL (1976) Digital coding of speech in sub-bands. *Bell Syst Tech J* 55:1069–1085.