

## MIT Open Access Articles

*Schema learning for the cocktail party problem*

The MIT Faculty has made this article openly available. **Please share** how this access benefits you. Your story matters.

**Citation:** Woods, Kevin J. P., and Josh H. McDermott. "Schema Learning for the Cocktail Party Problem." Proceedings of the National Academy of Sciences 115, no. 14 (March 21, 2018): E3313–E3322.

**As Published:** <http://dx.doi.org/10.1073/PNAS.1801614115>

**Publisher:** Proceedings of the National Academy of Sciences

**Persistent URL:** <http://hdl.handle.net/1721.1/119661>

**Version:** Final published version: final published article, as it appeared in a journal, conference proceedings, or other formally published context

**Terms of Use:** Article is made available in accordance with the publisher's policy and may be subject to US copyright law. Please refer to the publisher's site for terms of use.





# Schema learning for the cocktail party problem

Kevin J. P. Woods<sup>a,b</sup> and Josh H. McDermott<sup>a,b,1</sup>

<sup>a</sup>Department of Brain and Cognitive Sciences, Massachusetts Institute of Technology, Cambridge, MA 02139; and <sup>b</sup>Program in Speech and Hearing Bioscience and Technology, Division of Medical Sciences, Harvard University, Boston, MA 02115

Edited by Richard M. Shiffrin, Indiana University, Bloomington, IN, and approved February 21, 2018 (received for review February 2, 2018)

**The cocktail party problem requires listeners to infer individual sound sources from mixtures of sound. The problem can be solved only by leveraging regularities in natural sound sources, but little is known about how such regularities are internalized. We explored whether listeners learn source “schemas”—the abstract structure shared by different occurrences of the same type of sound source—and use them to infer sources from mixtures. We measured the ability of listeners to segregate mixtures of time-varying sources. In each experiment a subset of trials contained schema-based sources generated from a common template by transformations (transposition and time dilation) that introduced acoustic variation but preserved abstract structure. Across several tasks and classes of sound sources, schema-based sources consistently aided source separation, in some cases producing rapid improvements in performance over the first few exposures to a schema. Learning persisted across blocks that did not contain the learned schema, and listeners were able to learn and use multiple schemas simultaneously. No learning was evident when schema were presented in the task-irrelevant (i.e., distractor) source. However, learning from task-relevant stimuli showed signs of being implicit, in that listeners were no more likely to report that sources recurred in experiments containing schema-based sources than in control experiments containing no schema-based sources. The results implicate a mechanism for rapidly internalizing abstract sound structure, facilitating accurate perceptual organization of sound sources that recur in the environment.**

auditory scene analysis | perceptual learning | implicit learning | statistical learning

Sounds produced by different sources sum in the air before entering the ear, requiring the auditory system to infer sound sources of interest from a mixture (the “cocktail party problem”) (1–9). Because many different sets of source signals could generate an observed mixture, the problem is inherently ill-posed. In the real world, however, constraints on the generation of sound mean that assumptions can be made about which components of sound energy came from the same source, enabling us to correctly infer source structure much of the time. Understanding human listening abilities thus requires understanding these assumptions and how they are acquired.

Some of the assumptions guiding scene analysis may be rather general. For example, frequency components appearing at integer multiples of a common “fundamental frequency” are usually heard as arising from the same source (10–12), as are sounds that begin and end at the same time (13, 14) and sound patterns that repeat (15–17). These grouping cues are believed to reflect constraints on sound generation that are common across natural sources (1) and thus likely apply across a wide range of sounds and contexts.

Other cues to perceptual organization might apply only to particular contexts. Natural sources often produce sounds that are patterned consistently across occurrences (as in animal vocalizations, spoken words, or sung melodies), resulting in an abstract time-varying structure shared by a subset of sound events. Internalizing this recurring structure might be expected to aid scene analysis, but unlike more generic grouping cues, which could be internalized over evolution or by a general learning process operating on all auditory input, source-specific structure would have to be learned upon the appearance of a new sound source.

Although auditory memory has been argued to have lower capacity than visual memory (18), human listeners clearly acquire

rich knowledge of sound structure from listening. Many documented examples fall under the rubric of “statistical learning,” in which humans internalize aspects of the sound input distribution, such as transition probabilities between sound elements (19–21) or correlations between sound properties (22). Such learning is thought to be important for both speech (23) and music (24–26) perception. Specific recurring sound structures, typically noise samples, can also be learned (27–29). Such learning is apparently often implicit (24–26, 30).

The ability to learn the structure of sound sources suggests that such knowledge might be used for scene analysis, and source-specific structures used for this purpose are often termed “schemas” in the scene-analysis literature (1, 31–36). A role for learned schemas has been suggested by prior findings that listeners are better able to extract highly familiar voices (e.g., one’s spouse) (37), familiar languages (38), well-known melodies (39–42), and words (43). However, because these sources were already familiar to listeners before the experiments, the underlying learning process has remained opaque. Open issues include the rapidity with which schemas can be learned and used in scene analysis, the specificity of the learned representation, whether schemas can be learned in the presence of multiple sources, whether learning is dependent on attention to schema exemplars, and whether listeners must be aware of what they are learning. Also, because prior work has largely been confined to familiar structures in speech and music, it has been unclear if schema learning is a general phenomenon in audition.

The experiments presented here were designed to reveal the process of learning a new schema. Our approach was to have listeners perform source-separation tasks on synthetic stimuli that traversed a pattern over time and to test if performance improved for targets derived from a particular pattern (the schema) that appeared intermittently over the course of the experiment. We employed this general approach in three separate experimental

## Significance

**The “cocktail party problem” is encountered when sounds from different sources in the world mix in the air before arriving at the ear, requiring the brain to estimate individual sources from the received mixture. Sounds produced by a given source often exhibit consistencies in structure that might be useful for separating sources if they could be learned. Here we show that listeners rapidly learn the abstract structure shared by sounds from novel sources and use the learned structure to extract these sounds when they appear in mixtures. The involvement of learning and memory in our ability to hear one sound among many opens an avenue to understanding the role of statistical regularities in auditory scene analysis.**

Author contributions: K.J.P.W. and J.H.M. designed research; K.J.P.W. performed research; K.J.P.W. analyzed data; and K.J.P.W. and J.H.M. wrote the paper.

The authors declare no conflict of interest.

This article is a PNAS Direct Submission.

Published under the PNAS license.

<sup>1</sup>To whom correspondence should be addressed. Email: jhm@mit.edu.

This article contains supporting information online at [www.pnas.org/lookup/suppl/doi:10.1073/pnas.1801614115/-DCSupplemental](http://www.pnas.org/lookup/suppl/doi:10.1073/pnas.1801614115/-DCSupplemental).

Published online March 21, 2018.

paradigms with different types of stimuli that varied in complexity but never contained familiar source structure. In paradigm 1, listeners discriminated four-tone melodies presented concurrently with “distractor” tones, or comparable stimuli composed of noise bursts (Fig. 1). In paradigm 2, listeners had to attentively track one of two concurrent sources that changed stochastically over time in pitch and the first two formants (spectral peaks that determine vowel quality) (Figs. 2–4). In paradigm 3, pitch and formant contours were extracted from recorded speech and resynthesized to produce a stimulus that contained the pitch and formant contours of an actual speech utterance but that was not intelligible. Listeners heard a mixture of two such utterances followed by a probe utterance and were asked if the probe utterance was contained in the mixture (Fig. 5). Audio demonstrations of all stimuli can be found online at [mcdermottlab.mit.edu/schema\\_learning/](http://mcdermottlab.mit.edu/schema_learning/).

In each of these experimental paradigms, sources generated from a common schema recurred over the course of the testing session. These schema-based sources never appeared on consecutive trials and were transformed each time to avoid exact replication and to mimic the variation that occurs in real-world sound sources. We then compared performance for sources derived from a common schema with that for sources derived from schemas that did not recur during an experiment. The results show that rapidly acquired memories contribute substantially to source separation.

## Results

**Paradigm 1. Detection of Discrete-Tone Melodies.** To explore schema learning with a relatively simple stimulus, we presented listeners with a six-tone “mixture” composed of a four-note melody with two additional distractor tones, followed by a four-tone “probe” melody in isolation (Fig. 1A). Listeners were asked if the isolated probe melody was contained in the mixture that preceded it. The probe melody was always transposed away from the melody in the mixture by up to an octave, and listeners were told that this would be the case. The transposition required listeners to extract the structure of the melody and prevented them from performing the task based on glimpsed features (e.g., note fragments) of the mixture. On trials where the correct response was “no” (the probe was not contained in the mixture), the probe was altered by changing the middle two tones, with the first and last notes of the melody retaining the relative positions they had in the mixture. As a consequence, the task could not be performed based on these outer tones alone. The tone onsets and durations in the probe melody were unaltered on these “foil” trials so that the task also could not be performed by recognizing temporal patterns alone. Because we wanted to explore the learning of novel structure, melodies were not confined to a musical scale or metrical grid; pitch and timing values were drawn from continuous uniform distributions so that there was no conventional musical structure.

A schema-based melody appeared in the mixture on every other trial (Fig. 1B) and on half of those trials also appeared as the four-tone probe. Although the recurring schema could thus occur in isolation (as the probe), the alternating-trials design meant that a schema-based probe never immediately preceded a mixture containing that schema, preventing immediate priming. The non-schema-based trials for each participant consisted of trials drawn randomly from the schema-based sets for other participants (one from each of the other sets), so that schema- and non-schema-based stimuli were statistically identical when pooled across participants. As a consequence, any difference in performance between schema- and non-schema-based trials must reflect learning of the schema.

Because pilot experiments indicated that learning effects might be rapid, it seemed desirable to run large numbers of participants on relatively short experiments. The number of participants required was beyond our capacity to run in the laboratory, and so we instead recruited and ran participants online using Amazon’s

Mechanical Turk service. To mitigate concerns about sound quality, we administered a headphone-screening procedure to detect participants disregarding our instructions to wear headphones (44). Evidence that the quality of data obtained online can be comparable to that from the laboratory was obtained by comparing performance between online and in-laboratory participants, described below (experiment S2).

**Schema learning of melodies (experiment 1).** Listeners performed 100 trials of this task (taking ~10–15 min to complete). If exposure to a recurring melodic structure can help listeners detect it when it is embedded among distractors, then we might expect performance on schema-based trials to exceed that on non-schema-based trials over the course of the experiment.

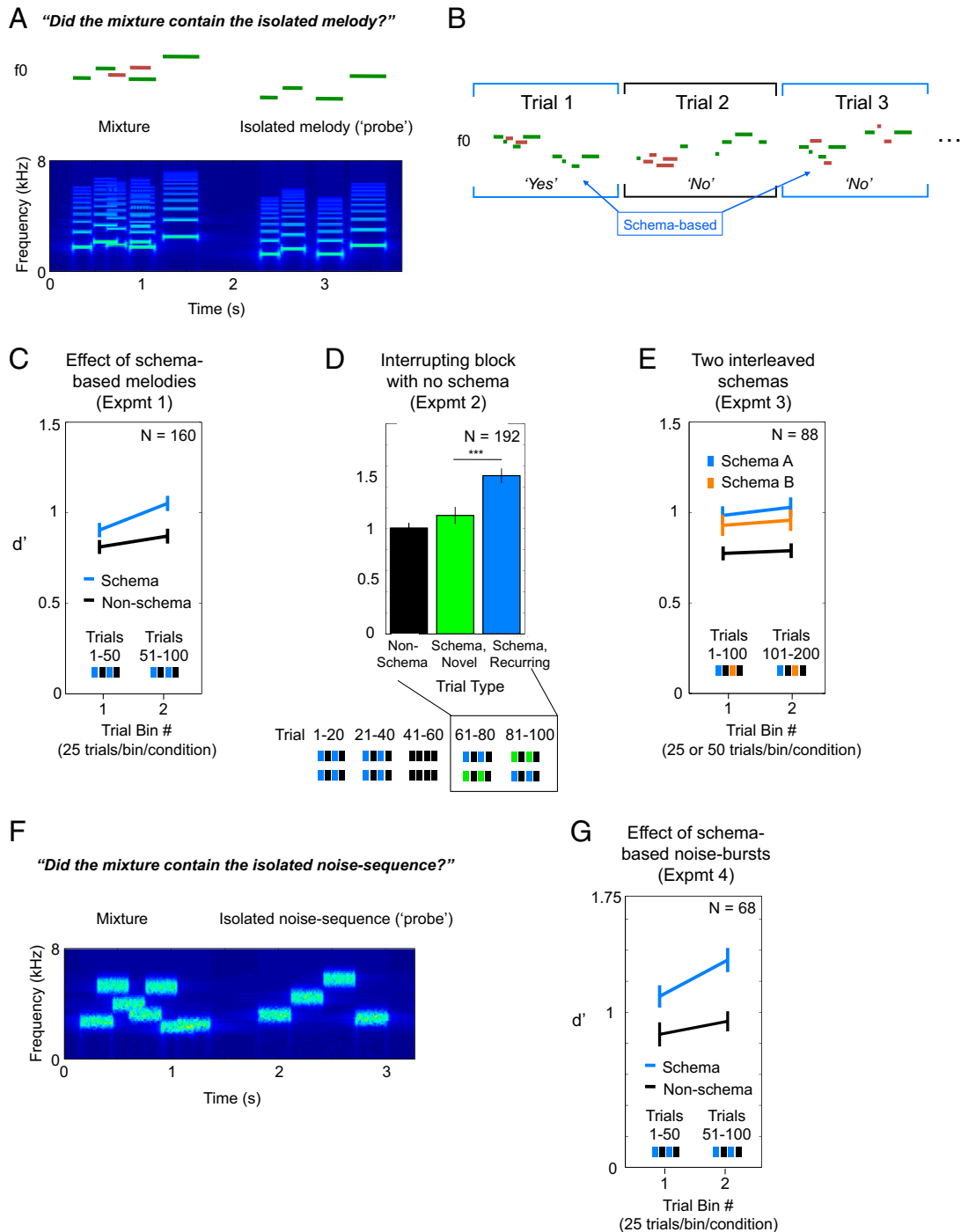
As shown in Fig. 1C, performance improved over time for both schema- and non-schema-based sources [ $F(1,159) = 5.06, P = 0.026$ ] but was better for schema- than for non-schema-based trials [ $F(1,159) = 7.69, P = 0.0062$ ]. Because the schema- and non-schema-based stimuli were statistically identical when pooled across participants, the performance benefit for schema-based trials indicates that participants learned and applied the structure of the schema. There was no interaction between trial type and time [ $F(1,159) = 0.97, P = 0.32$ ], perhaps because the learning of the schema was relatively rapid (we return to this issue in paradigm 2). However, the difference between schema and non-schema performance was significant in the second half of the session [ $t(159) = 2.61, P = 0.0098$ ] but not during the first half of the session [ $t(159) = 1.47, P = 0.14$ ]. The above-chance performance even in the absence of a schema indicates some ability to match the probe stimulus to the mixture despite the transposition. However, this ability is augmented by the acquired knowledge of the tones that are likely to belong together (the schema).

Because schema learning should, at a minimum, produce an improvement in performance late in the experiment, in most subsequent experiments we test for learning by comparing performance between conditions in the second half of trials within an experiment. Figures accordingly plot results binned into a small number of time bins, typically two (to maximize power).

**Schema learning is persistent (experiment 2).** To test whether knowledge of the schema could be retained over time, we conducted a second experiment in which exposure to a schema was followed by a middle block in which the schema-based melody was totally absent, which in turn was followed by blocks featuring the original schema or a new schema (Fig. 1D). The interrupting block contained 20 trials (~3 min). When melodies based on the original schema returned in the final block, they showed a performance benefit compared with the new schema [ $t(189) = 3.49, P < 0.001$ ]. This benefit suggests that effects of exposure early in the experiment persisted across the interrupting middle block.

**Multiple schemas can be learned simultaneously (experiment 3).** The persistence of a learned schema might allow multiple schemas to be learned and used concurrently. To examine this possibility, we conducted a third experiment in which two different schemas alternated, again interspersed with non-schema trials (each schema thus appeared on every fourth trial). The experiment was lengthened to 200 trials to present each schema 50 times, as before. As shown in Fig. 1E, the results suggest a learning effect for each schema similar to what we saw in the previous experiments [pooled schema-based trials vs. non-schema-based trials in second half of experiment;  $t(87) = 2.98, P < 0.005$ ]. These results suggest that multiple schemas can be learned and used at the same time.

**Schema learning occurs without isolated exposure to the schema (experiment S1).** To test whether the learning effects were dependent on exposure to the schema in isolation via its presence in the probe stimulus, we conducted an experiment which did not present isolated probes. Instead, listeners were presented with two mixtures and judged whether they contained the same melody. Performance was



**Fig. 1.** Schema learning in melody segregation (paradigm 1). (A) Schematic of the trial structure (*Upper*) and a spectrogram of a sample stimulus (*Lower*). A target melody (green line segments) was presented concurrently with two distractor notes (red line segments), followed by a probe melody (green line segments). Listeners judged whether the probe melody matched the target melody in the mixture. The probe melody was transposed up or down in pitch by a random amount. (B) Schematic of the basic experiment structure. On every other trial the target melody was generated from a common schema. On schema-based trials, the melody in the mixture was drawn from the schema 50% of the time, while the probe was always drawn from the schema. (C) Results of experiment 1: recognition of melodies amid distractor tones with and without schemas ( $n = 160$ ). Error bars throughout this figure denote the SEM. (D) Results of experiment 2: effect of an intervening trial block on learned schema ( $n = 192$ ). Listeners were exposed to a schema, then completed a block without the schema, and then completed two additional blocks, one containing the original schema and one containing a new schema. The order of the two blocks was counterbalanced across participants. (*Lower*) The two rows of the schematic depict the two possible block orders. (E) Results of experiment 3: effect of multiple interleaved schemas ( $n = 88$ ). Results are plotted separately for the two schemas used for each participant, resulting in 25 and 50 trials per bin for the schema and non-schema conditions, respectively. (F) Spectrogram of a sample stimulus from experiment 4. Stimulus and task were analogous to those of experiment 1, except that noise bursts were used instead of tones. (G) Results of experiment 4: recognition of noise-burst sequences amid distractor bursts, with and without schemas ( $n = 68$ ).

low overall for this experiment, presumably because there were twice as many opportunities to make streaming errors (Fig. S1). However, the schema learning effect persisted, with better performance for schema-based than for non-schema-based trials in the second block [ $t(39) = 3.10, P < 0.005$ ]. Listeners thus appear to be able to detect and learn recurring structure even when it does not occur in isolation.

**Schema learning occurs for atypical sound sources (experiment 4).** To test whether comparable phenomena would occur for sound sources that were even less typical of musical sources, we conducted an analogous experiment with sequences of noise bursts. Unlike the tones, the noise bursts were aperiodic and lacked a pitch in the usual sense but nonetheless instantiated patterns of frequency variation that were recognizable to human listeners (45). As shown in Fig. 1G, listeners again performed better for stimuli generated from a common schema [second block;  $t(67) = 3.56, P < 0.001$ ]. It thus appears that there is some generality to the ability to learn recurring patterns and use this knowledge to improve the extraction of those patterns from mixtures with other sounds.

**Paradigm 2. Attentive Tracking of Smooth Pitch-Formant Contours.** To further explore the generality of this phenomenon, we next turned to a task and stimulus we had originally developed to study auditory attentive tracking (46). Synthetic voices were generated that varied in several speech-relevant feature dimensions (pitch and the first two formants:  $f_0, F_1, F_2$ ) according to independent, randomly generated trajectories. Listeners were presented with mixtures of two such time-varying sources and were cued beforehand to attend to one of them (with the starting portion of one voice). We measured listeners' ability to track this cued voice by subsequently presenting them with the tail end of one of the voices; their task was to judge whether this probe segment belonged to the cued voice (Fig. 2A).

Trajectories in a mixture were required to cross each other at least once in each feature dimension, so listeners could not perform the task simply by attending to a high or low value of one of the features. Although the trajectories of each source were continuous, task performance is not critically dependent on complete continuity, as it is robust to the insertion of intermittent gaps in the stimuli (46). Moreover, despite the continuity, the task is effortful for human listeners, and success depends in part on accurate tracking of the cued voice as it varies throughout the mixture (46). One other difference between paradigms 1 and 2 was that the probe in paradigm 2 consisted only of the ending portion of a source (unlike paradigm 1, in which the probe had the same form as the target melody). As a result, listeners never experienced a source trajectory in the absence of a concurrent source, providing another test of whether schemas can be learned and used even when sources never occur in isolation.

Given that we ran experiments online to obtain sufficient sample sizes, it is natural to wonder whether data quality was comparable to that in experiments run in the laboratory. To address this issue, we compared performance on this attentive tracking task between online and in-laboratory participants (experiment S1). We chose to perform this comparison for the attentive tracking paradigm because it seemed most likely to suffer in a subject pool that was less motivated, as might be expected in an online environment. However, we found that performance was similar between online and in-laboratory participants once online participants were screened for headphone use (Fig. S2). This result gave us some confidence in our online testing procedures.

**Schema learning extends to pitch-formant contours (experiment 5).** To first test for the basic schema learning effect with this task, we ran an experiment in which the cued voice on every other trial was derived from a common schema trajectory (Fig. 2B). These schema-based sources were not exact replicas of each other but were related by time dilation and transposition, as might occur in natural sound sources, such as prosodic patterns in speech. Trials in which the target was not schema-based had targets drawn from

the sets of schema-based targets presented to other participants (Methods), so that when pooled across subjects the distribution of schema-based and non-schema-based targets was identical. To better explore the time course of any learning effect, we ran a longer experiment than we did for paradigm 1 (168 trials; ~35–40 min, which we divided into two time bins for analysis with maximum power but also plot in six bins to provide a sense of the dynamics over time).

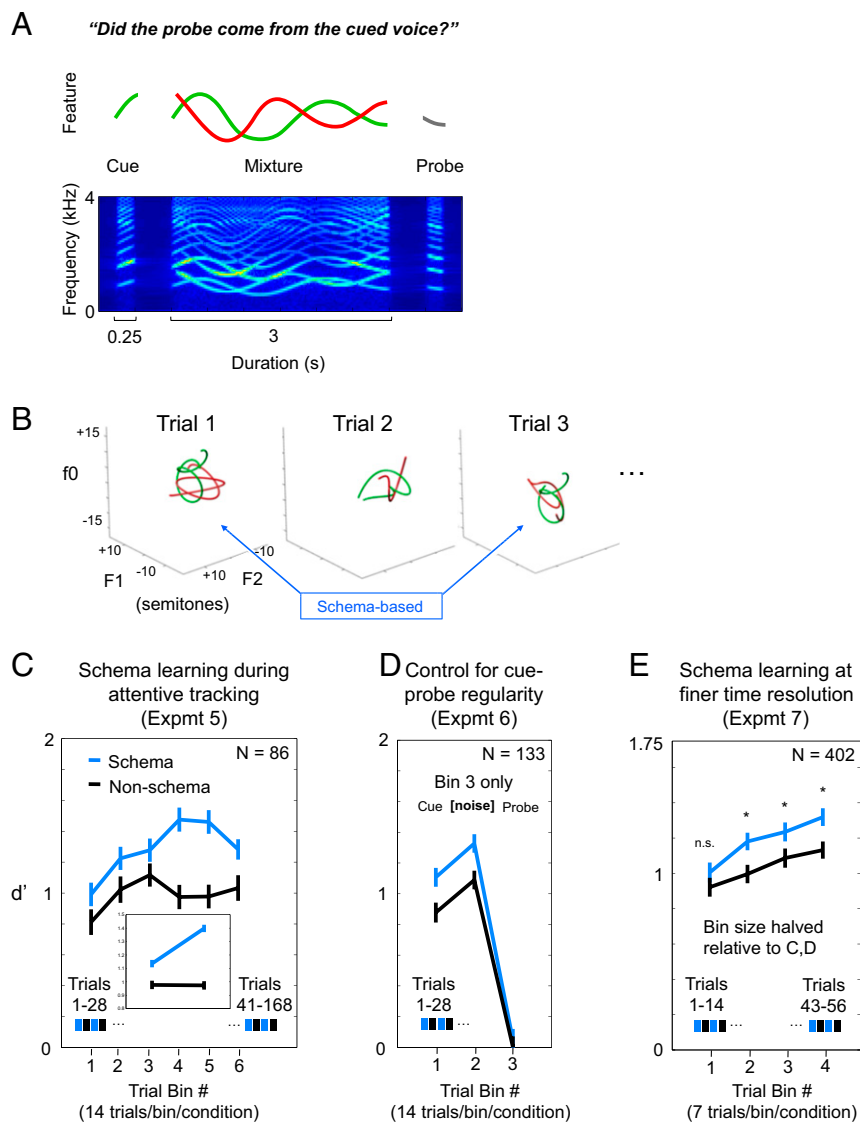
Performance over the course of the experiment is shown in Fig. 2C. Overall task performance was again significantly better for trials whose targets were based on a common schema [main effect of trial type,  $F(1,85) = 16.2, P = 0.0001$ ; repeated-measures ANOVA]. Moreover, although there was a general improvement over the course of the experiment [ $F(1,85) = 10.1, P = 0.002$ ], performance on schema-based target trials improved more than did performance for regular trials, yielding a significant interaction of trial type and time [ $F(1,85) = 12.5, P = 0.0007$ ; repeated-measures ANOVA], again driven by a performance difference in the second half of trials [ $t(85) = 5.14, P < 10^{-5}$ ]. These results suggest that performance can be facilitated by recurring structure even when sources vary in multiple dimensions and never appear in isolation.

**Learning effect is not explained by cues and probes (experiment 6).** One potential explanation for the learning effects in this task is that listeners learn something about the relationship between the cues and probes for the schema-based trajectories rather than from the trajectory itself. Although time dilation of schema-based trajectories resulted in the associated cues and probes not having a fully consistent relation to one another, we nonetheless guarded against any such strategy when generating stimuli by matching the distributions of distances between the cues and the correct and incorrect probes (Methods). To ensure that these measures indeed prevented listeners from performing the task with the cue and probe alone, we ran a second experiment in which the mixture on the last one-third of trials was replaced by noise (Fig. 2D). In a pilot experiment we found that when the relationship between the cues, correct probes, and incorrect probes was not controlled, performance remained significantly above chance during the noise block for both types of trials [non-schema trials:  $t(87) = 4.41, P < 0.000$ ; schema trials:  $t(87) = 3.22, P < 0.0018$ ], demonstrating the effectiveness of this control experiment (and the necessity of controlling the stimuli).

Replicating experiment 5, superior performance for schema-based trajectories was apparent over the first two-thirds of the experiment [ $F(1,132) = 8.35, P = 0.005$ ] (Fig. 2D). However, performance fell to chance levels once the mixtures were replaced by noise [ $t$  tests vs. chance: schema-based trajectories,  $t(132) = 0.71, P = 0.48$ ; non-schema-based trajectories,  $t(132) = 0.06, P = 0.96$ ], with no difference in performance between schema-based and non-schema-based trajectories [ $t(132) = 0.46, P = 0.65$ ]. It thus appears that listeners cannot perform the attentive tracking task based on the cue and probe alone, and that the benefit of schema-based trajectories is not due to learning cue–probe relationships for these trajectories.

**Rapid learning evident via crowdsourcing (experiment 7).** In experiments 5 and 6, the schema-based sources seemed to have elicited different levels of performance from the outset of the experiment (Fig. 2C and D). Because the balancing of the stimulus sets was intended to prevent an intrinsic difference in difficulty between conditions, we considered the possibility that learning might be occurring on a fast timescale. To test this, we pooled the data from experiments 4 and 5 (these experiments were identical for the first 56 trials) and examined performance over smaller bins of seven trials per bin rather than 14. A power analysis indicated that additional participants would be required to discover possible effects at this timescale, and so an additional 183 participants were run on a shorter, 56-trial version of the attentive tracking paradigm (experiment 7).

With a resolution of seven trials per bin, it is evident that performance at the outset of the experiment did not differ

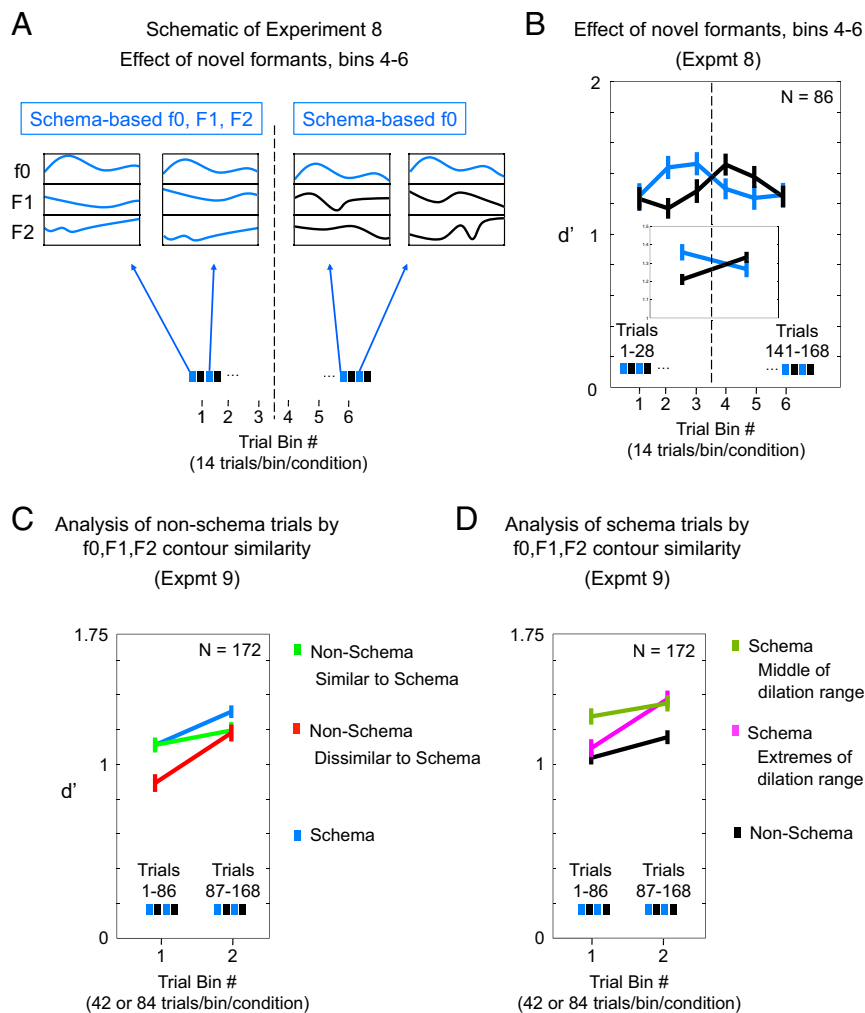


**Fig. 2.** Schema learning in attentive tracking of synthetic voices (paradigm 2). (A) Schematic of the trial structure (*Upper*) and spectrogram of an example stimulus (*Lower*). A target voice (green curve) was presented concurrently with a distractor voice (red curve). Both voices varied smoothly but stochastically over time in three feature dimensions:  $f_0$ ,  $F_1$ , and  $F_2$  (the fundamental frequency and first two formants; for clarity the schematic only shows variation in a single dimension). Voices in a mixture were constrained to cross at least once in each dimension. Listeners were cued beforehand with the initial portion of the target voice. Following the mixture, listeners were presented with a probe stimulus that was the ending portion of one of the voices and judged whether this probe came from the target. (B) Schematic of the experiment structure. On every other trial the target voice was generated from a common schema. Voices are depicted in three dimensions.  $f_0$ ,  $F_1$ , and  $F_2$  are plotted in semitones relative to 200, 500, and 1,500 Hz, respectively. (C) Results of experiment 5: effect of schemas on attentive tracking ( $n = 86$ ). The *Inset* denotes results with trials binned into 42 trials per condition to maximize power for an interaction test (reported in text). Error bars throughout this figure denote the SEM. (D) Results of experiment 6: a control experiment to ensure listeners could not perform the task with cues and probes alone ( $n = 146$ ). In the last one-third of trials, the voice mixture was replaced with noise. (E) Schema learning on a finer time scale ( $n = 402$ ). Data from the first 56 trials of experiments 5 and 6 were combined with new data from experiment 7 and replotted with seven trials per bin. The finer binning reveals similar performance at the experiment's outset, as expected. n.s., not significant. \* $P < 0.05$ .

between conditions [first time bin:  $t(401) = 0.92$ ,  $P = 0.36$ ] (Fig. 2E) but that performance differences emerge quickly (significant differences for all other time bins:  $P < 0.05$  in all cases). The ability to observe this rapid learning was facilitated by the fact that the experiments were run online, which allowed us to efficiently test a relatively large number of listeners ( $n = 402$ ). This observation provides some confirmation that the stimuli in the schema-based and non-schema-based conditions do not differ in their intrinsic difficulty; it is only the presence of other schema-based stimuli that boosts performance. The results are suggestive of a learning effect occurring over relatively small numbers of exposures.

**Schema learning need not be explicit.** The presence of a learning effect raises the question of whether participants are aware of what they are learning. To address this issue, after finishing the task, participants were asked if they had noticed repetition in the cued voice. The proportion of "yes" responses from experiment 5 (the longest experiment run with this paradigm) is shown in Fig. 3 along with responses from two control experiments: one that contained no schema-based sources (control experiment 1) and another in which every cued voice was schema-based (control experiment 2, in which we expected participants to notice the recurring source structure). Participants did not report repetition in experiment 5 any more often than in the control experiment with no schema-based sources





**Fig. 4.** Dependence of schema benefit on multiple dimensions and similarity to the schema. (A) Schematic of the structure of experiment 8. On each trial, listeners were cued to a target voice, heard a target-distractor voice pair, and judged if a subsequent probe was from the end of the target or the distractor (paradigm 2). Schema-based trials alternated with non-schema-based trials, but the formant trajectories on schema-based trials were randomized halfway through the experiment. (B) Results of experiment 8: effect of multiple dimensions on schema learning ( $n = 86$ ). The *Inset* denotes results with trials binned into 42 trials per condition to maximize power. Error bars throughout this figure denote the SEM. (C) Results of subdividing non-schema trials from experiment 9 ( $n = 146$ ). Performance was computed separately for non-schema trials whose feature trajectories were most and least correlated with those of the schema. (D) Results of subdividing schema trials from experiment 9. Performance was computed separately for schema trials in the middle and extremes of the dilation/transposition range.

non-schema-based trials were statistically identical when pooled across participants.

As shown in Fig. 5C, a benefit was rapidly obtained from the recurring schema, with performance on schema-based trials again exceeding non-schema-based trials during the second half of the experiment [ $t(92) = 2.94, P = 0.0041$ ]. These results demonstrate rapid schema learning for natural feature trajectories that have more complex generative constraints than the stimuli used in paradigms 1 and 2, raising the possibility that fast and flexible schema learning could help us hear behaviorally relevant sources in the real world.

**Schema-based distractors (experiment 11).** It seemed of interest to test effects on performance when schema-based sources appeared as the distractor instead of the target. However, paradigms 1 and 2 did not provide a clear means to study this: In paradigm 1 the distractors were not of the same form as the target melodies (being a pair of two tones rather than a four-note sequence), while in paradigm 2 the need to control cue-probe relationships made it methodologically challenging to implement experiments with similar distractors. The mixture-probe task of paradigm

3 was well suited to address this question, so we conducted an additional experiment in which the nontarget utterance in the mixture (on alternate trials) was generated from a common schema (Fig. 5D). This experiment was run for 200 trials rather than 100 trials as in experiment 7, providing the listener with considerable exposure to the schema by the end of the experiment. The results (Fig. 5E) show that the schema-based distractor nonetheless had no detectable effect on performance [ $F(1,201) = 1.85, P = 0.18$ ; no significant difference for any time bin], providing evidence that a recurring schema is less likely to be internalized if it does not occur in the attended source.

## Discussion

Sources in auditory scenes often produce sound with some degree of consistency. We explored the conditions in which this consistency might be learned and used to guide scene analysis. We tested if listeners would obtain a source-separation benefit from the recurrence of a source under transformations such as transposition and time dilation, which produce acoustically distinct variants that share abstract structure. Such a benefit would





As in some instances of visual implicit learning (49), learning nonetheless appears to be somewhat limited to task-relevant stimuli. We thus demonstrate that source structure can be learned amid concurrent sources but perhaps only when attention is directed to it.

**Relation to Effects of Short-Term Source Repetition.** Another relevant line of previous work involves effects of sources that exhibit regular or repeating structure. For example, cyclically repeating patterns in ambiguous streams of tones are known to group over time to form a single auditory stream (16). Repetition also causes random noise sources to be segregated from mixtures (17). These phenomena are distinct from those that we studied here in that the recurring structure is exact, occurs on a short timescale, causes the repeating elements to group together, and shows no evidence of learning (i.e., retention over periods where the source is not present). That said, one can envision a continuum between the conditions of these previous studies (back-to-back and exact repetition) and those of the present study (abstract recurrence across intervening stimuli), and it remains to be seen whether there is any relation between the underlying mechanisms.

**What is Learned?** Our listeners evidently learned something about the recurring structure in each experiment. Because sources were transposed and time-dilated/compressed over the course of learning, the recurring schemas were not individuated by particular feature values. Our results suggest that listeners instead learned something about the way the source's features changed over time. Experiment 8 demonstrated that the learned schema can incorporate variation in formants as well as pitch, since the performance benefit for schema-based stimuli was eliminated when formant trajectories were randomized in schema-based trials. Experiment 9 showed that the learned schema provided a benefit to non-schema targets that were sufficiently similar to the schema (provided the similarity metric included both pitch and formants) and an added benefit to schema targets in the middle of the range of possible time-dilated variants. Overall, the results indicate that learning can occur over a range of task-relevant features and that the effect of the schema is graded, providing the greatest benefit to stimuli most similar to the canonical schema.

The recurring structures in our experiments were abstract, since the schema always appeared transposed or dilated/compressed to varying degrees (transformations inspired by the variation that occurs in speech and music). It would be interesting to further explore the nature of the learned representation by testing the transfer of learning across different transformations (e.g., time reversal or rotation of trajectories in feature-space) and to explore limits to the types of abstract structure that can be learned (e.g., by exposing listeners to different types of source transformations during learning). There are also likely limits to the contexts in which they can be utilized. For example, listeners are known to have difficulty detecting even a highly familiar melody if it is perfectly interleaved with a set of distractor tones (31, 39). Understanding how schemas interact with other constraints on source separation is thus another important topic for future research.

**Primitive vs. Schema-Based Scene Analysis.** Schema-based scene analysis in audition has historically been contrasted with "primitive" scene analysis, in which putatively universal grouping cues are mediated by processes that do not require attention or memory (1, 8, 35, 50). For example, sequential sounds that are similar (e.g., in pitch or timbre) often arise from the same source in the world and tend to group perceptually over time (7, 35, 51). However, because schema-based scene analysis has not been studied extensively in the laboratory, little is known about the underlying mechanisms, and the extent to which they are distinct from those of

primitive scene analysis has been unclear. The methodology introduced here should enable future progress on these issues.

It is possible that the schemas that are learned in our experiments affect perception in much the same way as putatively primitive grouping cues (e.g., pitch differences between talkers). This notion could be tested by comparing the neural or behavioral consequences of schema-driven segregation with those of segregation via other cues (e.g., pitch differences). For instance, it could be diagnostic to examine whether a learned schema affects the ability to discriminate the temporal relationship between elements of the schema and another concurrent source, which is often taken as a signature consequence of streaming (43, 52, 53).

The effect of the learned schema may also be to alter the interaction of streaming and attention. It could be that a learned schema makes it easier to attend to a source conforming to the schema, explaining the better performance on our tracking task, for instance. Alternatively, if memory is complementary to attention, then schema learning might serve to reduce the attentional resources that would otherwise be required to segregate a recurring source from others. These possibilities could be disentangled by measuring attentional selection [for instance, by asking listeners to detect perturbations to sources (46)] before and after schema learning.

However, should the tasks we used even be considered to involve streaming? All the stimuli involve discriminating sound sources embedded in mixtures with other sounds, but the stimuli were relatively short. As such, they are distinct from the long sequences of alternating tones commonly used to study streaming (35, 51). Such stimuli notably exhibit a "buildup of streaming" in which the likelihood of hearing two streams increases with time (3, 53, 54). Although the stimuli we used do not evoke this particular phenomenon, they nonetheless require sound energy to be grouped over time. As such, we conceive them as involving streaming in a more general sense of the term and view them as useful for understanding real-world scenarios in which sources do not repeat cyclically ad infinitum.

The rapidity of the learning effects shown here also raise the possibility that learning could influence all aspects of scene analysis, even those that are quite general in their applicability. Even evidence that newborns exhibit aspects of similarity-based streaming (55, 56) is consistent with learning from early experience. The difference between primitive and schema-based processes might thus be better described in terms of the scale and scope of learning: Primitive scene analysis could effectively be mediated by schema that are very general and that can be applied indiscriminately.

**Schema Learning May Be Ubiquitous in Audition.** In real-world auditory scenes, sources are sometimes unfamiliar, and recurring structure may occur only intermittently and concurrent with other sounds. Our results demonstrate that the auditory system can rapidly learn to utilize abstract source structure even in such challenging conditions. The robustness of this learning could allow schema-based scene analysis to occur across a much wider range of scenarios than previously imagined.

## Materials and Methods

All experiments were approved by the Committee on the Use of Humans as Experimental Subjects at the Massachusetts Institute of Technology. On the initial page of the online task, participants read a disclaimer (as required by the MIT Committee for the Use of Humans as Experimental Subjects) and consented to participation in the experiment. Participants in the in-laboratory condition in experiment S1 signed a form indicating their consent to participate. Methods are described in [SI Materials and Methods](#).

**ACKNOWLEDGMENTS.** We thank the members of the J.H.M. laboratory for helpful comments on the manuscript. This work was supported by a McDonnell Foundation Scholar Award (to J.H.M.), National Science Foundation Grant BCS-1454094, NIH Grant 1R01DC014739-01A1, and NIH Training Grant T32DC000038.

1. Bregman AS (1990) *Auditory Scene Analysis: The Perceptual Organization of Sound* (MIT Press, Cambridge, MA).
2. Bronkhorst AW (2000) The cocktail party phenomenon: A review of research on speech intelligibility in multiple-talker conditions. *Acta Acust United Acust* 86: 117–128.
3. Carlyon RP (2004) How the brain separates sounds. *Trends Cogn Sci* 8:465–471.
4. 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.
5. McDermott JH (2009) The cocktail party problem. *Curr Biol* 19:R1024–R1027.
6. Shinn-Cunningham BG (2008) Object-based auditory and visual attention. *Trends Cogn Sci* 12:182–186.
7. Shamma SA, Micheyl C (2010) Behind the scenes of auditory perception. *Curr Opin Neurobiol* 20:361–366.
8. Snyder JS, Gregg MK, Weintraub DM, Alain C (2012) Attention, awareness, and the perception of auditory scenes. *Front Psychol* 3:15.
9. Middlebrooks JC, Simon JZ, Popper AN, Fay RR (2017) *The Auditory System at the Cocktail Party* (Springer, New York).
10. Moore BCJ, Glasberg BR, Peters RW (1986) Thresholds for hearing mistuned partials as separate tones in harmonic complexes. *J Acoust Soc Am* 80:479–483.
11. Hartmann WM, McAdams S, Smith BK (1990) Hearing a mistuned harmonic in an otherwise periodic complex tone. *J Acoust Soc Am* 88:1712–1724.
12. Micheyl C, Oxenham AJ (2010) Pitch, harmonicity and concurrent sound segregation: Psychoacoustical and neurophysiological findings. *Hear Res* 266:36–51.
13. Darwin CJ (1981) Perceptual grouping of speech components differing in fundamental frequency and onset-time. *Q J Exp Psychol Sect A* 33:185–207.
14. 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.
15. 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.
16. Bendixen A, Denham SL, Gyimesi K, Winkler I (2010) Regular patterns stabilize auditory streams. *J Acoust Soc Am* 128:3658–3666.
17. McDermott JH, Wroblewski D, Oxenham AJ (2011) Recovering sound sources from embedded repetition. *Proc Natl Acad Sci USA* 108:1188–1193.
18. Cohen MA, Horowitz TS, Wolfe JM (2009) Auditory recognition memory is inferior to visual recognition memory. *Proc Natl Acad Sci USA* 106:6008–6010.
19. Saffran JR, Johnson EK, Aslin RN, Newport EL (1999) Statistical learning of tone sequences by human infants and adults. *Cognition* 70:27–52.
20. Creel SC, Newport EL, Aslin RN (2004) Distant melodies: Statistical learning of non-adjacent dependencies in tone sequences. *J Exp Psychol Learn Mem Cogn* 30: 1119–1130.
21. Aslin RN, Newport EL (2012) Statistical learning: From acquiring specific items to forming general rules. *Curr Dir Psychol Sci* 21:170–176.
22. Stip CE, Rogers TT, Kluender KR (2010) Rapid efficient coding of correlated complex acoustic properties. *Proc Natl Acad Sci USA* 107:21914–21919.
23. Saffran JR, Newport EL, Aslin RN (1996) Word segmentation: The role of distributional cues. *J Mem Lang* 35:606–621.
24. Tillmann B, Bharucha JJ, Bigand E (2000) Implicit learning of tonality: A self-organizing approach. *Psychol Rev* 107:885–913.
25. Pearce MT, Ruiz MH, Kapasi S, Wiggins GA, Bhattacharya J (2010) Unsupervised statistical learning underpins computational, behavioural, and neural manifestations of musical expectation. *Neuroimage* 50:302–313.
26. Rohrmeier M, Rebuschat P (2012) Implicit learning and acquisition of music. *Top Cogn Sci* 4:525–553.
27. Kaernbach C (2004) The memory of noise. *Exp Psychol* 51:240–248.
28. Agus TR, Thorpe SJ, Pressnitzer D (2010) Rapid formation of robust auditory memories: Insights from noise. *Neuron* 66:610–618.
29. Leek MR, Watson CS (1988) Auditory perceptual learning of tonal patterns. *Percept Psychophys* 43:389–394.
30. Perruchet P, Pacton S (2006) Implicit learning and statistical learning: One phenomenon, two approaches. *Trends Cogn Sci* 10:233–238.
31. Bey C, McAdams S (2002) Schema-based processing in auditory scene analysis. *Percept Psychophys* 64:844–854.
32. Haykin S, Chen Z (2005) The cocktail party problem. *Neural Comput* 17:1875–1902.
33. Darwin CJ (2008) Listening to speech in the presence of other sounds. *Philos Trans R Soc Lond B Biol Sci* 363:1011–1021.
34. Alain C, Bernstein LJ (2008) From sounds to meaning: The role of attention during auditory scene analysis. *Curr Opin Otolaryngol Head Neck Surg* 16:485–489.
35. Moore BCJ, Gockel HE (2012) Properties of auditory stream formation. *Philos Trans R Soc Lond B Biol Sci* 367:919–931.
36. Bronkhorst AW (2015) The cocktail-party problem revisited: Early processing and selection of multi-talker speech. *Atten Percept Psychophys* 77:1465–1487.
37. Johnsrude IS, et al. (2013) Swinging at a cocktail party: Voice familiarity aids speech perception in the presence of a competing voice. *Psychol Sci* 24:1995–2004.
38. Cooke M, Garcia Lecumberri ML, Barker J (2008) The foreign language cocktail party problem: Energetic and informational masking effects in non-native speech perception. *J Acoust Soc Am* 123:414–427.
39. Dowling WJ (1973) The perception of interleaved melodies. *Cognit Psychol* 5:322–337.
40. Dowling WJ, Lung KM-T, Herrbold S (1987) Aiming attention in pitch and time in the perception of interleaved melodies. *Percept Psychophys* 41:642–656.
41. Devergie A, Grimault N, Tillmann B, Berthommier F (2010) Effect of rhythmic attention on the segregation of interleaved melodies. *J Acoust Soc Am* 128:EL1–EL7.
42. Szalárdy O, et al. (2014) The effects of rhythm and melody on auditory stream segregation. *J Acoust Soc Am* 135:1392–1405.
43. Billig AJ, Davis MH, Deeks JM, Monstrey J, Carlyon RP (2013) Lexical influences on auditory streaming. *Curr Biol* 23:1585–1589.
44. Woods KJP, Siegel MH, Traer J, McDermott JH (2017) Headphone screening to facilitate web-based auditory experiments. *Atten Percept Psychophys* 79:2064–2072.
45. McDermott JH, Lehr AJ, Oxenham AJ (2008) Is relative pitch specific to pitch? *Psychol Sci* 19:1263–1271.
46. Woods KJP, McDermott JH (2015) Attentive tracking of sound sources. *Curr Biol* 25: 2238–2246.
47. Nygaard LC, Pisoni DB (1996) Learning voices. *J Acoust Soc Am* 99:2589–2603.
48. Reinisch E, Holt LL (2014) Lexically guided phonetic retuning of foreign-accented speech and its generalization. *J Exp Psychol Hum Percept Perform* 40:539–555.
49. Jiang Y, Chun MM (2001) Selective attention modulates implicit learning. *Q J Exp Psychol A* 54:1105–1124.
50. Fritz JB, Elhilali M, David SV, Shamma SA (2007) Auditory attention—Focusing the searchlight on sound. *Curr Opin Neurobiol* 17:437–455.
51. Moore BCJ, Gockel H (2002) Factors influencing sequential stream segregation. *Acta Acust United Acust* 88:320–333.
52. Micheyl C, Hunter C, Oxenham AJ (2010) Auditory stream segregation and the perception of across-frequency synchrony. *J Exp Psychol Hum Percept Perform* 36: 1029–1039.
53. Thompson SK, Carlyon RP, Cusack R (2011) An objective measurement of the build-up of auditory streaming and of its modulation by attention. *J Exp Psychol Hum Percept Perform* 37:1253–1262.
54. Roberts B, Glasberg BR, Moore BCJ (2008) Effects of the build-up and resetting of auditory stream segregation on temporal discrimination. *J Exp Psychol Hum Percept Perform* 34:992–1006.
55. McAdams S, Bertoncini J (1997) Organization and discrimination of repeating sound sequences by newborn infants. *J Acoust Soc Am* 102:2945–2953.
56. Winkler I, et al. (2003) Newborn infants can organize the auditory world. *Proc Natl Acad Sci USA* 100:11812–11815.
57. Peer E, Brandimarte L, Samat S, Acquisti A (2017) Beyond the Turk: Alternative platforms for crowdsourcing behavioral research. *J Exp Soc Psychol* 70:153–163.
58. Peer E, Vosgerau J, Acquisti A (2014) Reputation as a sufficient condition for data quality on Amazon mechanical Turk. *Behav Res Methods* 46:1023–1031.
59. Dowling WJ, Fujitani DS (1971) Contour, interval, and pitch recognition in memory for melodies. *J Acoust Soc Am* 49:524–531.
60. Klatt DH (1980) Software for a cascade/parallel formant synthesizer. *J Acoust Soc Am* 67:971–995.
61. Garofolo JS, Lamel LF, Fisher WM, Fiscus JG, Pallett DS (1993) DARPA TIMIT acoustic-phonetic continuous speech corpus CD-ROM. NIST speech disc 1-1.1. *NASA STIRecon Tech Rep N* 93:27403.
62. Kawahara H, et al. (2008) Tandem-STRAIGHT: A temporally stable power spectral representation for periodic signals and applications to interference-free spectrum, F0, and aperiodicity estimation. *2008 IEEE International Conference on Acoustics, Speech and Signal Processing (IEEE, Piscataway, NJ)*, pp 3933–3936.
63. Mustafa K, Bruce IC (2006) Robust formant tracking for continuous speech with speaker variability. *IEEE Trans Audio Speech Lang Process* 14:435–444.
64. Buhrmester M, Kwang T, Gosling SD (2011) Amazon’s mechanical Turk: A new source of inexpensive, yet high-quality, data? *Perspect Psychol Sci* 6:3–5.
65. Crump MJC, McDonnell JV, Gureckis TM (2013) Evaluating Amazon’s mechanical Turk as a tool for experimental behavioral research. *PLoS One* 8:e57410.
66. Paolacci G, Chandler J (2014) Inside the Turk understanding mechanical Turk as a participant pool. *Curr Dir Psychol Sci* 23:184–188.