

Infant biases for detecting speech in complex scenes

Christina M. Vanden Bosch der Nederlanden^a, Athena Vouloumanos^b

^aThe University of Western Ontario, ^bNew York University

Author Note

^a The Brain and Mind Institute, Department of Psychology, The University of Western Ontario; ^b Department of Psychology, New York University

Correspondence should be addressed to Christina Vanden Bosch der Nederlanden, The Brain and Mind Institute, Department of Psychology, The University of Western Ontario, 1151 Richmond St, London, ON N6A 3K7. cdernede@uwo.ca
This research was supported by the Eunice Kennedy Shriver National Institute of Child Health & Human Development of the National Institutes of Health under Award Number R01HD072018 awarded to AV. Special thanks to all of the members of the NYU Infant Cognition and Communication Lab, and especially all of the families who participated in this study.

The authors declare no conflicts of interest.

Data Availability Statement: All de-identified data are available upon request.

Abstract

How do infants learn the sounds of their native language when there are many simultaneous sounds competing for their attention? Adults and children detect when speech sounds change in complex scenes better than when other sounds change. We examined whether infants have similar biases to detect changes to human speech better than non-speech sounds including musical instruments, water, and animal calls in complex auditory scenes. We used a change deafness paradigm to examine whether 5-month-olds' change detection is biased toward certain sounds within high-level categories (e.g., biological, or generated by humans) or whether change detection depends on low-level salient physical features such that detection is better for sounds with more distinct acoustic properties, such as water. In Experiment 1, 5-month-olds showed some evidence of detecting changes to speech and music sounds better than no change trials. In Experiment 2, when speech and music were compared separately with animal and water sounds, infants detected when speech and water changed, but not when music changed across scenes. Infants' change detection is both biased for certain sound categories, as they detected small changes to speech better than other comparable sounds, and affected by the size of the acoustic change, similar to young infants' attentional priorities in complex visual scenes. By 5 months, infants show some preferential processing of speech changes in complex auditory environments which could help bootstrap the language learning process.

Keywords: Auditory scene analysis; change deafness, infant change detection, attentional bias, speech, category knowledge, perceptual salience

1 Introduction

2 A critical early step for building infants' language system is to selectively attend
3 to speech sounds even when other sounds are presented simultaneously. Infants' early
4 and powerful speech perception abilities have been explored through decades of
5 innovative laboratory studies, including investigations of speech preference (e.g.,
6 Vouloumanos & Werker, 2004; 2007), discrimination (Kuhl, 2004; Werker & Tees, 1984;
7 Werker, Yeung, Yoshida, 2012), segmentation (Saffran, Aslin, & Newport, 1996;
8 Jusczyk, Houston, & Newsome, 1999; Thiessen, Hill, & Saffran, 2005), and
9 categorization (Waxman, Markow, 1995; Werker et al., 1998; Swingley, 2009; Lim,
10 Lacerda, & Holt, 2015). Laboratory studies, however, are typically conducted with
11 speech presented in quiet, even sound-attenuated, conditions. In our daily interactions,
12 human speech is often heard overlapping with other sounds, like dogs, cats, airplanes,
13 construction, wind, or rain, which also compete for our attention (Bregman, 1990).
14 Prioritizing speech sounds in our environment is therefore critical for processing and
15 learning about speech (Vouloumanos & Waxman, 2014). We examined whether
16 processing biases for detecting speech in complex auditory scenes are evident in
17 infancy, when they could provide a potential mechanism for facilitating language
18 learning.

19 Adults and children show biases toward social stimuli in their environment. From
20 birth, infants are biased to attend toward face or face-like stimuli (e.g., Farroni et al.,
21 2005; Simion et al., 2001) and infants' bias for attending to faces in dynamic visual
22 scenes increases with age during their first year (Frank, Vul, & Johnson, 2009). In the
23 auditory domain, newborns listen longer to speech sounds compared to carefully

1 matched non-speech stimuli in silent, controlled environments (Vouloumanos & Werker,
2 2004; 2007). Infants' listening bias for speech becomes more specific over the first 3
3 months, with infants preferring to listen to speech compared to monkey calls, and other
4 human-produced non-speech sounds, like throat clearing or laughing (Shultz &
5 Vouloumanos, 2010). Adults and children as young as 6 years old are better at
6 encoding and detecting changes involving the human voice in complex auditory scenes
7 (Vanden Bosch der Nederlanden, Snyder, & Hannon, 2016). However, very few studies
8 have examined infants' speech perception in these noisier, ecological conditions more
9 typical of human experience (cf., Newman & Jusczyk, 1996; Newman et al., 2006). If
10 processing biases for speech in complex scenes were evident during infants' first year,
11 they could provide a potential mechanism for facilitating language learning. That is,
12 becoming an expert at speech perception may require infants to selectively attend to
13 human communicative sounds in auditory scenes that have multiple sounds competing
14 for their attention.

15 Two key factors affecting how humans attend to sounds in our environments are
16 salience and semantics. In visual scene processing, 3-month-olds attend more to low-
17 level, salient physical features than high-level (but less salient), semantically meaningful
18 combinations of features specific to faces compared to 6- and 9-month-olds (Frank, Vul,
19 & Johnson, 2009), suggesting it may take time for infants to fully develop a bias for
20 faces *qua* faces over salient changes within a visual scene. Adults' and children's
21 detection of changing sounds in complex auditory scenes is driven by these same two
22 factors as their processing of complex visual scenes. Although large acoustic changes
23 are easier to detect than small acoustic changes, mature listeners rely more on their

semantic knowledge of real-world sounds than acoustics to detect changes between scenes (Gregg, Irsik, & Snyder, 2014; Vanden Bosch der Nederlanden, Snyder, & Hannon, 2016). Further, adults are better at encoding and detecting changes to human vocal sounds, and speech in particular, compared to human non-speech sounds, musical instruments, animal calls, and environmental sounds, regardless of the size of the acoustic change (Vanden Bosch der Nederlanden et al., 2018). Experienced listeners use their knowledge of real-world sound categories to preferentially attend to speech sounds in complex scenes (Gregg & Samuel, 2009), but it is unclear whether relatively inexperienced language learning infants would have similar processing biases for detecting changes to speech in complex scenes.

We examined whether infants have biases for processing speech in complex scenes. Specifically, we contrasted infants' detection of changes in speech sounds with their detection of changes in various non-speech sounds in complex auditory scenes. We adapted a research design from the auditory change deafness literature (Vitevitch, 2003; Eramudugolla et al., 2005; see Gregg & Snyder, 2011) in which listeners hear two auditory "scenes" with multiple real-world sounds presented at the same time and must determine whether the sounds in each scene were the same or different. We compared change detection for speech to non-speech sound categories—musical instruments, animal calls, and water sounds—to test whether infants' ability to notice a change is influenced by high-level specific sound categories or low-level physical properties like the size of the acoustic change. We measured the size of the acoustic change in a two-dimensional space by calculating the Euclidian distance between changing sounds within a trial (see Acoustic analyses in Method). For example, a change between a man

speaking and a cello melody might be very similar in pitch (frequency) and would result in a small perceptual change. In contrast, changing from a woman speaking to a tuba melody would be a big change in pitch and timbre, which may be easier to detect.

In Experiment 1, we compared speech, musical instruments, and water sounds. Acoustic analyses showed that speech and musical instruments have similar acoustic features, allowing us to test whether—given equivalent acoustic distance—infants detect changes to speech better than other sounds. Water is more acoustically distant from both speech and musical instruments, allowing us to test how larger acoustic distances influenced infant detection. In Experiment 2, we further examined how sound categories and physical salience affect infant change detection by presenting speech and music to different infants and contrasting them with changes in animal sounds which have similar acoustic features, and water which has larger acoustic distance compared to the three other sounds. This study allowed us to compare the detection of music and speech changes in isolation from one another, while also examining the effect of physical salience and sound category on change detection in infancy. Across these two experiments, different non-speech sound categories allowed us to contrast whether infants' change detection was better predicted by physical salience (large acoustic changes), as for young infants' visual scene perception, or by a bias for speech, as for adults' and children's auditory scene perception.

We predicted that if infants' auditory biases take time to develop as in visual face processing studies, 5-month-olds may detect large acoustic changes (e.g., to water) better than changes to other sounds, including speech. At the same time, if infants already have biases for processing speech over other acoustically similar sound

categories, they may preferentially detect changes to speech compared to music and animal calls (Vouloumanos & Werker, 2007). Better speech detection in complex scenes would suggest that infants' speech bias could help them isolate speech sounds and learn language in real-world auditory settings.

Method

Experiment 1

Participants. Participants were 30 healthy, full-term 5-month-old infants (13 females; *M* age: 5 months, 7 days; range 4 months, 15 days to 6 months, 1 day)¹. Participants were recruited from maternity wards at local hospitals, and were reported as being healthy and having normal hearing at the time of testing. All parents reported that their children were learning English, and 19 parents reported additional languages for their infants including Spanish (6), French (3), Italian (3), Hebrew (2), Yiddish (1), Cantonese and Mandarin (1), Vietnamese (1), German (1), and Polish (1). An additional 13 infants (9 females) were excluded from analyses because of infant fussiness (3), experimenter error (5), looking at the screen for the maximum trial length on 5 or more trials (2), medical history (2), and parental interference (1). Parents gave informed

¹ Previous studies examining change deafness and different trial types have found η_p^2 ranging from .30 (Vanden Bosch der Nederlanden et al., 2018) to .79 (Vanden Bosch der Nederlanden et al., 2015). We used the a priori G*Power (Faul et al., 2009) estimate of within-subjects main effects for ANOVA, which yielded a sample size of 10 adult participants for 2 groups (set1 and set 2) and 4 measurements (4 trial types). Next, we estimated the effect size from a previous infant study looking at attention to different sound types (Shultz & Vouloumanos, 2010, Exp 1) to be .13 and this yielded a minimum sample size of 18 participants. We aimed for a conservative sample size of 30, because of the variability in infant responses and because this was the first study of infant change deafness and their sensitivity had not yet been established.

consent on behalf of their infants and received a certificate and small toys or t-shirts as gifts for their participation. All procedures were approved by the IRB at New York University (IRB-FY2016-81, Divergent biases for conspecifics as early markers for autism spectrum disorders).

Stimuli. Sounds were drawn from three categories: human speech, musical instruments, and water sounds (see acoustic analyses in Results). Two sounds from each category (woman’s voice, man’s voice; cello, tuba; bubbling water, draining water) were used to create 18 trials, 9 trials for each set with one sound from each category in each set (Set 1: woman’s voice, cello, bubbling water; Set 2: man’s voice, tuba, draining water). The speech tokens were repeated English CV syllables with a speech-like contour, i.e., “de de de de” or “ma ma ma ma” spoken to the prosody of “I love that one!” in a positive adult-directed manner. Trials were composed of 2-object scenes. A 2-object “scene” consisted of 2 overlapping sounds 1000 ms in length with simultaneous onsets and offsets. These scenes were concatenated with 500 ms silent inter-stimulus intervals to make 43-second long sound trials.

Design. To examine infant processing of complex auditory scenes, we modified the one-shot procedure used in change deafness studies (see e.g., Snyder, Gregg, Weintraub, & Alain, 2012) and adapted the alternating procedure common in visual change blindness studies (e.g., Rensink, O’Regan, & Clark, 1997). This alternating paradigm switches between two static images with a brief blank image placed between them, with the length of time taken to detect the change indexing change blindness (longer reaction times indicate worse performance). This paradigm bears striking resemblance to infant looking time methods using stimulus-alternation preference

procedures (Cowan, Suomi, & Morse, 1982; Best & Jones, 1998), that take advantage of infants' preference for novelty, with longer looking times to trials that alternate between two different visual displays or phonemes compared to looking at visual displays or phonemes from a single category. In our auditory paradigm, infants heard either alternating 2-sound scenes (see Figure 1), called change trials, or non-alternating scenes, called no change trials. In a no change trial, the same 2-object scene was repeated for the entire length of the trial. In a change trial, the 2-object scene alternated with a different 2-object scene in which one of the sounds was the same and one of the sounds changed. For example, in a speech change trial, the woman's voice + draining water scene would alternate with the cello + draining water scene (see Figure 1A-C). Each infant heard 9 trials in total: 3 no change trials and 6 change trials: 2 speech change trials, 2 music change trials, and 2 water change trials (see Appendix for specific sound combinations). As described in the stimulus section, two sets of stimuli were created with different exemplars of each sound category. Half the infants were tested with Set 1 and half the infants were tested with Set 2. Some trials involved a change in the speech sound, others involved changes in different non-speech sounds, e.g. music. Infants' looking time (compared with change trials) indexed whether or not they noticed the change between complex auditory scenes. Comparing changes for different sound categories allows us to examine whether infants' change detection is based on the size of the acoustic change (e.g., better for water) or on the sound category that is changing (better for speech, musical instruments).

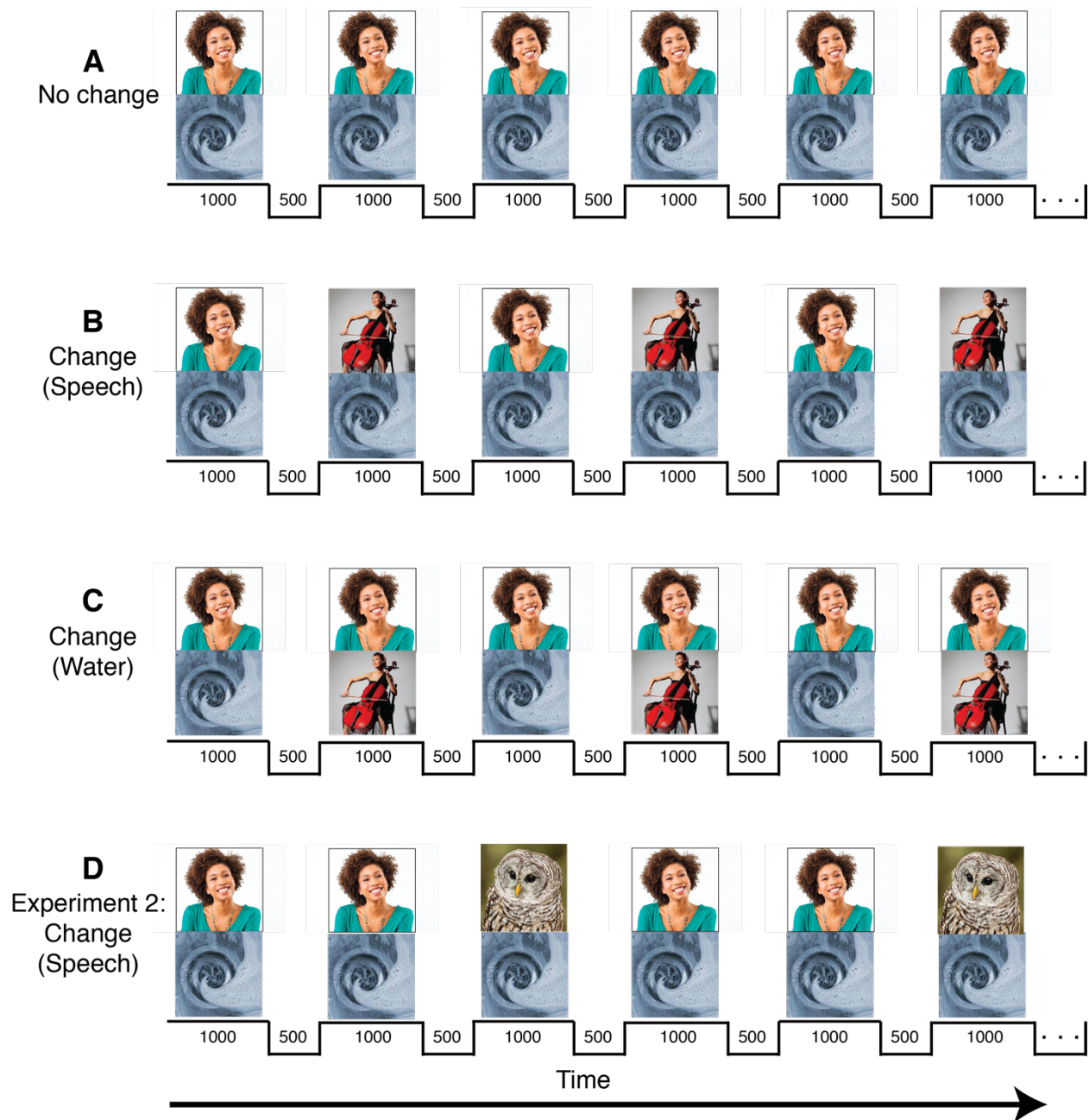


Figure 1. Example of no change and change trials for Experiment 1 (A-C) and Experiment 2 (A and D), which had a same-same-different structure instead of simple alternations between scenes.

Procedure. Infants were tested in a sound-attenuated room using an infant-controlled sequential preferential looking procedure (e.g., Cooper & Aslin, 1994, Vouloumanos & Werker, 2004), run in Habit 2.1.25 (Oakes, Sperka, & Cantrell, 2015). In this procedure, infants controlled the onset and offset of each trial by looking at or away from a central monitor. Infants sat on a parent's lap 35" (89 cm) in front of a 30" (76.25 cm) computer monitor. Parents wore headphones playing a sound medley to mask the experimental sounds. At the start of the experiment, infants' attention was drawn to the monitor by a colorful expanding and contracting circle. Once infants fixated on the monitor, a stationary black and white checkerboard appeared in tandem with one set of sounds presented at a mean amplitude of 60 dB (± 5 dB). An experimenter, blind to the stimulus condition, recorded infants' fixations to and away from the monitor during the study. Trials ended if infants looked away for longer than 2 seconds or if the maximum trial length was reached (43 seconds). A new trial began once infants looked at the expanding and contracting circle on the center monitor.

Infant looking time was coded online and verified offline by coding frame-by-frame looking to the monitor using 30 frame-per-second videos. The more precise offline looking times were used in the analyses. Average looking times for each of the four trial types were calculated for each participant: no change, speech change, music change, water change. Trials with missing looking time values were not include in the final looking time average for each participant. All figures are reported using within-subjects error bars given the repeated-measures design of our study (Cousineau, 2005).

Results and Discussion

Looking times for the sound dropping out of the first scene² were submitted to a repeated measures Analysis of Variance (ANOVA) with no change, speech change, music change, and water change trials as levels of the within-subjects factor sound and stimulus set (set 1, set 2) as a between-subject factor. There was no main effect ($p = .623$) or interaction ($p = .459$) with stimulus set. There was a main effect of trial type $F(3,84) = 4.425$, $p = .013$, $\eta_p^2 = .136$. Planned comparisons first examined whether infants detected the three sound category changes relative to no change trials: Infants looked longer at speech change ($M = 16.10$, $SD = 1.47$; $p = .019$) and music change trials ($M = 19.46$, $SD = 2.14$; $p = .005$) than no change trials, but water change trials ($M = 14.54$, $SD = 1.73$) did not differ from no change trials ($M = 13.24$, $SD = 1.17$; $p = .355$). Music change and speech change trials did not differ from one another ($p = .102$, see Figure 2). Water change trials did not differ from speech change trials ($p = .376$), but looking time was significantly higher for music than water change trials ($p = .041$).

Together these results suggest that infants looked longer at music and speech changes than no change trials. That is, when infants listened to multiple sounds presented simultaneously in a scene, they noticed changes to music and human speech equally well. At the same time, detection of changes in music, but not in human speech,

² Each change trial inherently involves a sound dropping out of one scene and a sound replacing it in the next scene. As such, change trials can be categorized by either the dropping or replacing sound. For instance, in Figure 1B, speech drops out of the first scene and is replaced by the cello. This change trial could be categorized as either a speech (dropping sound) or musical instrument (replacing sound) change trial. All analyses were completed based on both dropping and replacing sound groupings. Replacing sounds resulted in a marginal main effect of trial type, $F(3,84) = 2.682$, $p = .067$, $\eta_p^2 = .087$. The effect size for dropping sounds from the first scene is reported in the main text.

was better than detection of changes in water sounds suggesting some attentional priorities for music in auditory scene processing. Attentional priorities may be important for guiding infants' and children's attention toward relevant sounds in our environment and for scaffolding learning in the real world (e.g., Frank, Vul, & Johnson, 2009; Scerif, 2010).

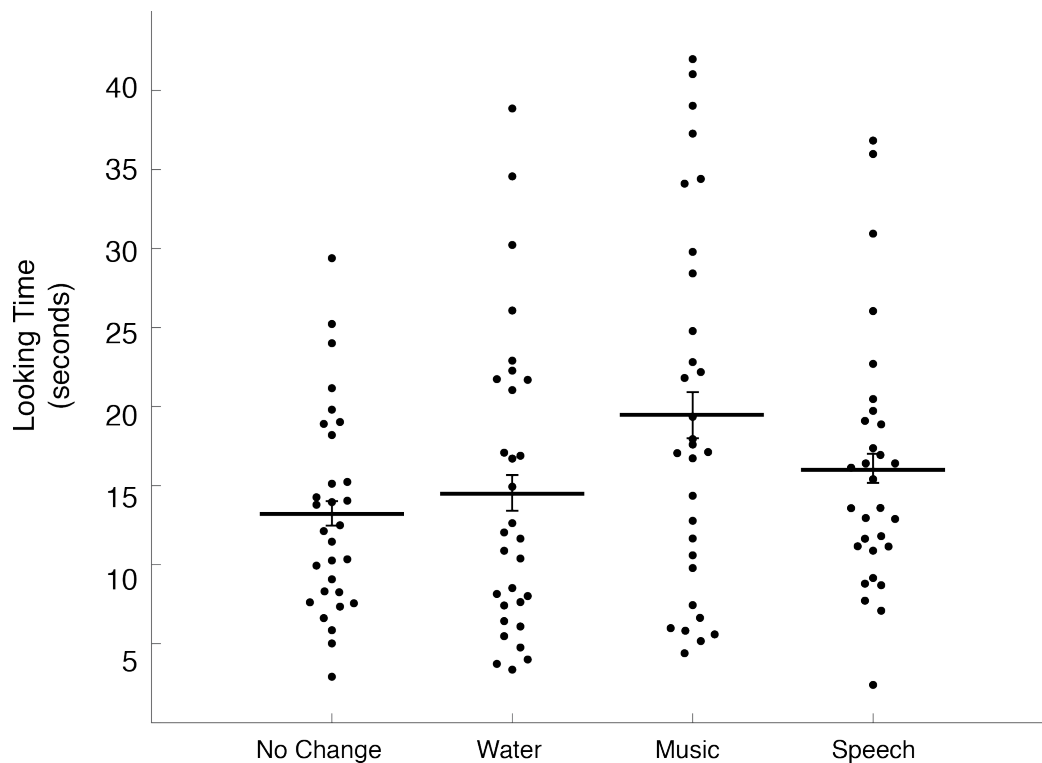


Figure 2. Infants looked longer at music and speech change trials than no change trials; water change trials were not different from no change trials. Speech change trials also did not differ from water change trials. Error bars are within-subjects standard error.

However, these results are not entirely straight-forward as speech change trials did not differ significantly from water change trials and water change trials did not differ from no change trials. The design in of Experiment 1 made it difficult to tease apart

looking time for speech and music separately as they were often heard together and, in some trials, music was heard when speech changed and vice-versa. Further, the alternating design suggested two possibilities for how to analyze change trials, grouping results either by the sound dropping or replacing in each scene (see Footnote 2).

Experiment 2

The results of Experiment 1 suggest that music and speech were both salient changes to infants, but the interpretation of this effect is tempered by the lack of a difference between speech and water change trials, and in turn, the lack of difference between water and no change trials. To further examine change detection of different sound categories, in Experiment 2, we separated music and speech into different conditions and adopted a same-same-different alternating design. The same-same-different alternating paradigm has been used in previous visual change deafness studies and more clearly highlights the stable sound as it occurs twice in every triplet. It is important to note that speech and musical instruments are well-matched in their acoustic features and so differences across conditions would not be due to pitch, harmonicity, or the size of the acoustic change. We compared speech and music to animal sounds because their acoustic features and the size of the acoustic change are comparable (see Table 1).

Participants. Participants were 36 healthy, full term 5-month-old infants (10 females; *M* age: 5 months, 12 days, range 4 months, 15 days to 6 months, 0 days).³ None of the participants had a history of ear infections or had a cold/ear infection at the time of the study. Parents reported that all children were learning English, and 29 parents reported additional languages for their infants including with Spanish (17), Hebrew (2), Russian (2), Italian (2), German (2), Farsi (1), Fulani (1), Tigrinya (1), and Hindi (1). An additional 6 infants were excluded from the current analyses because of equipment failure such that they did not hear any sound during the study (4), or because of fussiness (2). Infants were randomly assigned to either the music condition (17 infants, 4 females, mean age 5 months, 12 days) or the speech condition (19 babies, 6 females, mean age 5 months, 12 days). Parents gave informed consent on behalf of their infants and received a certificate and small toys or t-shirts as gifts for their participation. All procedures were approved by the New York University (IRB-FY2016-81, Divergent biases for conspecifics as early markers for autism spectrum disorders).

Stimuli. Sounds were drawn from four sound categories: human speech, musical instruments, water sounds, and animal sounds. Speech and music were presented in separate conditions. For speech, two sounds of each type (woman's voice, man's voice; bubbling water, draining water; owl hoot, sheep baa) were used to create 18 trials, 9 trials for each set, with one sound from each category in each set (Set 1: man's voice, owl hoot, draining water; Set 2: woman's voice, sheep baa, bubbling water). For music, two sounds of each type (cello, tuba; bubbling water, draining water; sheep baa, owl

³ Sample sizes were collected based off of the η_p^2 obtained from the Experiment 1 $\eta_p^2=.136$, using G*Power (Faul et al., 2009) to calculate that a total of 18 participants were needed per condition.

hoot) were used to create 18 trials, 9 trials for each set with one sound from each category in each set (Set 1: cello, owl hoot, draining water; Set 2: tuba, sheep baa, bubbling water). All trials were composed of 2-object scenes created in the same manner as Experiment 1, with 1000 ms overlapping sounds with simultaneous onsets and offsets with 500-ms silent intervals.

Design. Experiment 2 used the same design as Experiment 1. However, unlike Experiment 1, change scenes had a same-same-different sequence using 2-object scenes, with 10 repetitions of that sequence. No change trials consisted of a same-same sequence (Figure 1A). All trials were 45 seconds in length. For example, in a speech change trial, the woman's voice + draining water scene would be presented twice, followed by the owl + draining water scene (see Figure 1D). Each infant heard 9 trials in total: 3 no change trials and 6 change trials consisting of 2 speech or music change trials, 2 water change trials, and 2 animal change trials. Infants were pseudo-randomly assigned to the music or speech conditions. About half of the infants in each study were assigned to one of two stimulus sets (Speech: Set 1 = 9 infants, Set 2 = 10 infants; Music: Set 1 = 8 infants, Set 2 = 9 infants).

Procedure. Identical to Experiment 1.

Results and Discussion

To confirm that infants were equally attentive during both sound conditions, total looking time across all trial types was compared between the speech and music conditions in a one-way ANOVA with condition as a between-subjects factor. There was

1 no difference between total looking time for the music (55.82 s) and speech (60.53 s)
 2 conditions, $F(1, 34) = .406$, $p = .528$, $\eta_p^2 = .012$.

3 We analyzed whether looking time to change trials differed as a result of sound
 4 category and/or acoustic salience in two separate ANOVAs for the speech and music
 5 conditions. Looking time was analyzed with trial type (no change, water change, animal
 6 change, speech or music change) as the within-subjects variable, and stimulus set (Set
 7 1, Set 2) as a between subject variable.

8 **Speech.** There was a main effect of trial type $F(3,51) = 3.990$, $p = .024$, $\eta_p^2 =$
 9 $.190$, but no main effect for stimulus set, $F(1,17) = .596$, $p = .451$, $\eta_p^2 = .034$, or
 10 interaction, $F(3,51) = .137$, $p = .888$, $\eta_p^2 = .008$. Planned comparisons first examined
 11 whether infants detected the three sound changes relative to no change trials: speech
 12 change ($M = 18.86$, $SD = 2.22$; $p = .005$) and water change trials ($M = 17.43$, $SD = 2.39$;
 13 $p = .033$) differed from no change trials ($M = 11.58$, $SD = 1.17$), but animal change trials
 14 did not ($M = 12.65$, $SD = 1.42$; $p = .430$). Three planned comparisons examined how
 15 looking time for change trials differed from each other to examine how sound category
 16 and the magnitude of acoustic change influenced looking time. Speech change trials
 17 elicited longer looking times than animal change trials ($p = .045$), but did not differ from
 18 water change trials ($p = .614$) and animal change trials did not differ from water change
 19 trials ($p = .132$, see Figure 3). Thus, infants looked longer at speech trials than no
 20 change and animal change trials, but they looked equally at speech and water change
 21 trials.

22

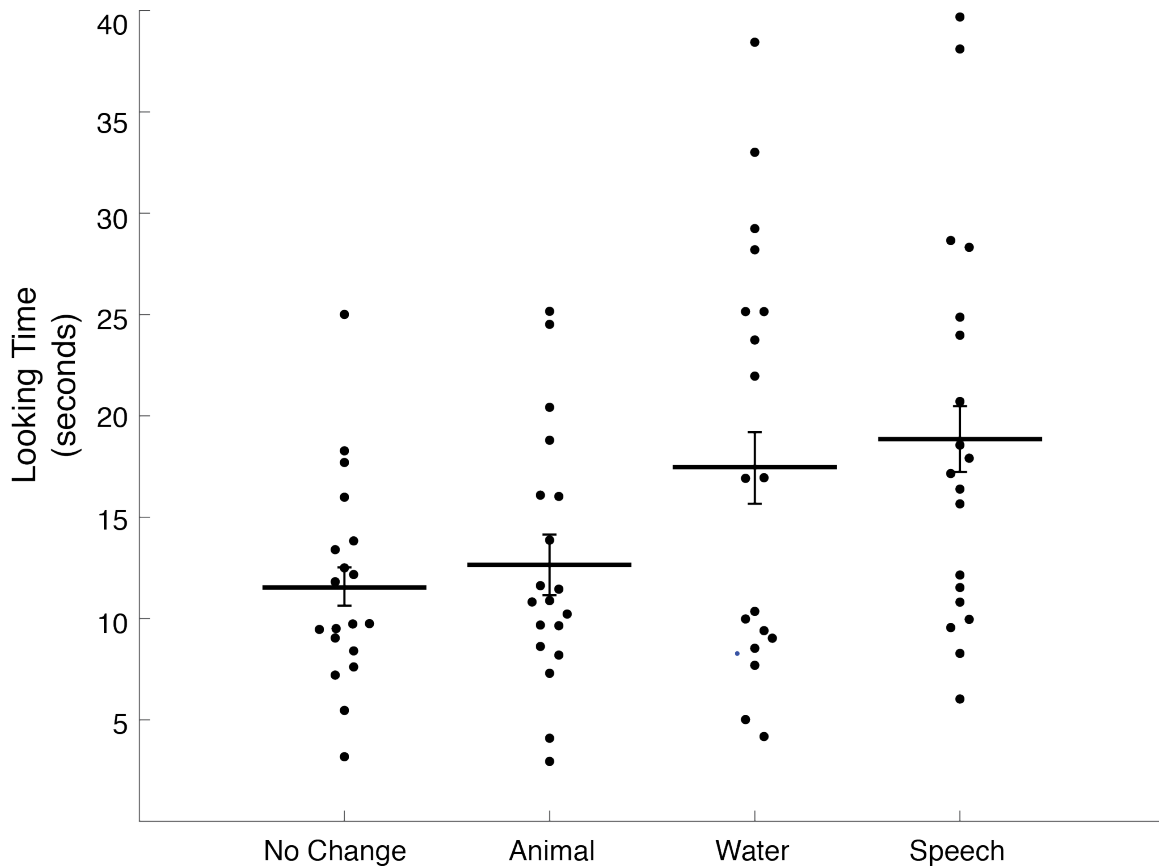


Figure 3. Mean looking times to no change, animal, water, and speech trials in Experiment 2. Error bars are within-subjects standard error.

Music. There were no main effects of trial type, $F(3,45) = 1.651$, $p = .191$, $\eta_p^2 = .099$ or set, $F(1,15) = 1.250$, $p = .281$, $\eta_p^2 = .077$, and no interaction ($p = .155$). These findings suggest that there was no difference between music change trials ($M = 15.28$, $SD = 1.32$) and no change trials ($M = 11.35$, $SD = 1.40$) or any other change trials ($M_{\text{animal}} = 13.15$, $SD = 1.53$; $M_{\text{water}} = 16.04$, $SD = 1.78$) in Experiment 2 (see Figure 4).

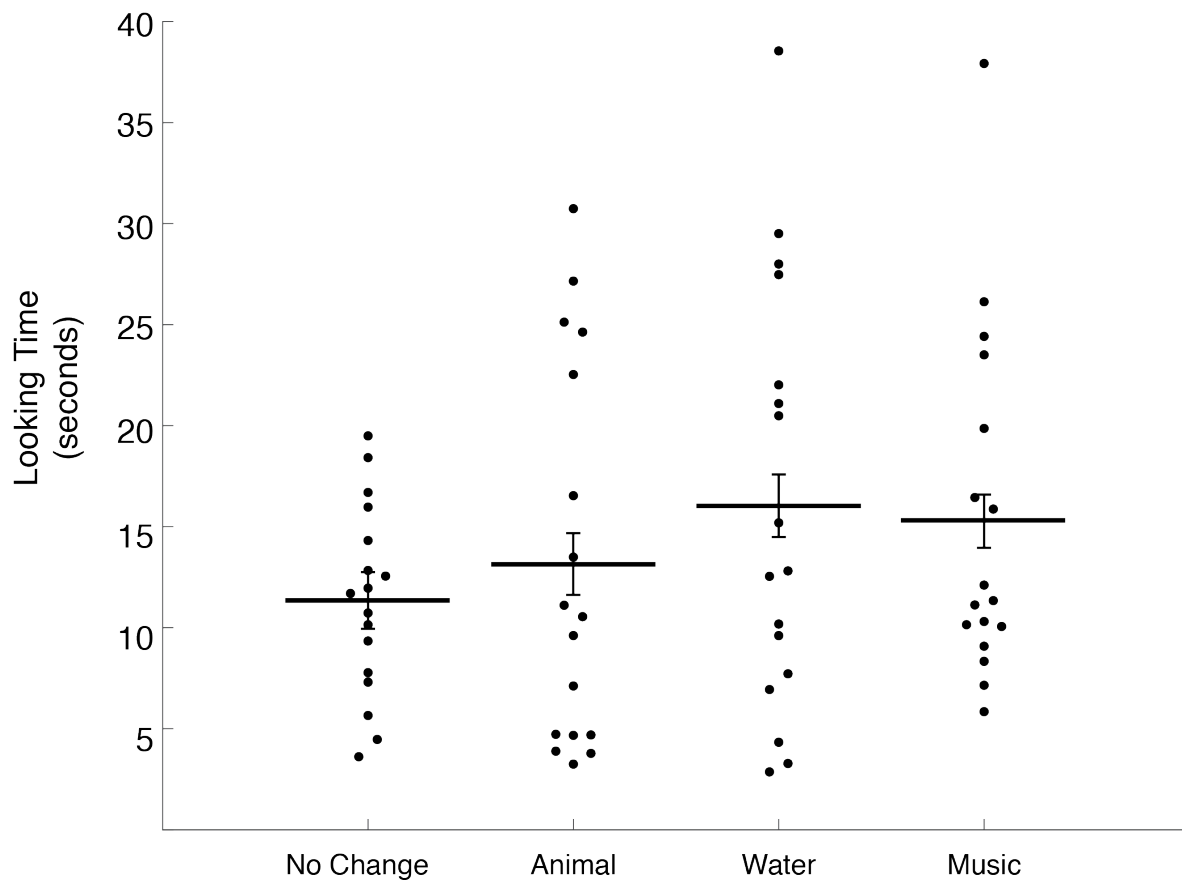


Figure 4. Mean looking times to no change, animal, water, and music trials in Experiment 2. Error bars are within-subjects standard error.

Together these findings suggest that infants attend to changes to speech sounds more than changes to music sounds in complex scenes composed of multiple sounds, despite both sound categories having similar acoustic properties. As in Experiment 1, speech change trials did not differ from water change trials. We further analyzed whether infants' looking time was influenced by the magnitude of the change for different sound types.

Acoustic analyses. The magnitude of the acoustic change between the dropping and replacing auditory object was calculated for all change trials using pitch and harmonicity as the two dimensions of this Euclidian space. These dimensions were selected based on previous change deafness studies (Gregg & Samuel, 2009, Vanden Bosch der Nederlanden, Hannon, & Snyder, 2016) and based on previous studies of environmental sound similarity (Gygi, Kidd, & Watson, 2007). Speech and musical instrument sounds were nearly identical in terms of their average pitch and harmonicity, animal sounds were very similar in their harmonicity, and water sounds differed from most other sounds in pitch and harmonicity (see Table 1). Euclidian distances were similar for speech, music, and animal change trials, but water change trials, because of their different acoustic features, resulted in a larger Euclidian distance than other sound categories (see Table 2). By including water sounds, we were able to look at how the size of the acoustic change influenced looking time.

Table 1. Acoustic characteristics (means) of each sound category: pitch, log of pitch, and harmonicity.

Stimuli	Pitch (Hz)	Log of Pitch	Harmonicity (dB)
Speech	162.92	2.19	14.63
Musical Instrument	153.87	2.19	13.42
Water	486.12	2.66	4.84
Animal	305.81	2.46	16.09

- 1 Table 2. Euclidian distances (mean of arbitrary units) measure the size of the acoustic
- 2 change for each sound category in each experiment.

Change Type	Expt 1: Speech and Music	Expt 2: Speech Only	Expt 2: Music Only
Speech	7.37	5.92	–
Musical Instrument	6.76	–	5.65
Water	9.20	10.54	9.93
Animal	–	6.64	6.98

- 3
- 4 **Euclidian distance.** Euclidian distance, the size of the acoustic change within a
- 5 change trial, may provide one possible explanation for the longer than expected looking
- 6 times to water sounds. The average size of the change was larger for water change
- 7 trials compared to the other change types (see Table 1 and Table 2). To test this
- 8 explanation, we examined whether infant looking times correlated with Euclidian
- 9 distance. There was a significant correlation between Euclidian distance and
- 10 participants' looking time for the speech condition, $r(18) = .585$, $p = .011$, but not the
- 11 music condition, $r(18) = .367$, $p = .135$. This suggests that infants were sensitive to the
- 12 size of the acoustic change in the speech condition, but not the music condition. Infants'
- 13 looking time did not correlate with Euclidian distance for the combined speech and
- 14 music scenes from Experiment 1, $r(18) = .338$, $p = .170$. Looking time may only have
- 15 correlated with the speech condition because this had the largest spread in Euclidian
- 16 distances compared to the other two conditions (see Table 2), suggesting that infants

may be sensitive to the size of the change, but require a larger spread in acoustic change magnitude to show this sensitivity.

The findings from Experiment 2 suggest that infants are biased toward listening to human speech in complex scenes, but they are also sensitive to acoustic salience, as they detected water changes in both conditions. Infants do not appear to show any such bias for music changes. Infants' equal detection of music changes and speech changes in Experiment 1 may have been due to speech also being played during the music change trials which heightened infants' processing in these scenes. By presenting speech and music sounds separately, we could examine the independent contribution of these different sound categories to differences in looking time in Experiment 2.

General Discussion

How do infants learn language in the real world when, so often, many different sounds are competing for infants' attention? We examined whether infants' early emerging listening biases for speech may be evident in scene processing, giving them an advantage in detecting and processing speech from noisy environments. In two experiments, 5-month-olds looked longer when speech sounds changed across auditory scenes compared to when acoustically comparable non-speech sounds changed. In Experiment 1, infants looked longer for changes in music and speech compared to no change trials. In Experiment 2, when infants heard either changes to speech or to music in separate conditions, infants looked longer for speech changes, but not for music changes compared to no change trials. This suggests that the longer looking times for music and speech sounds observed in Experiment 1 may have been due to both

sounds being presented within the same scene. At the same time, in both experiments, infants looked equally at changes in speech and at changes in water sounds, which had large Euclidian distances, suggesting that acoustic distance also influences how infants detect changes in sounds. This study provides the first evidence that infants preferentially process changes to speech in complex acoustic scenes with multiple overlapping sounds, as well as being sensitive to low-level physical properties such as acoustic distance.

Given comparable acoustic distance between sounds within sound categories, infants looked longer to speech changes than changes in music, and animal sounds. This is consistent with previous work with children and adults showing that listeners of all ages detect change by relying heavily on their knowledge of sound categories (Vanden Bosch der Nederlanden, Hannon, Snyder, 2016; Vanden Bosch der Nederlanden et al., 2018). Five-month-olds appear to use auditory listening strategies similar to older children and adults, biasing attention toward changes in socially relevant categories, even when those acoustic changes were small. Listening preferentially to the sounds that are most communicatively relevant can make the task of processing sounds within a complex scene much more efficient. Instead of having to encode and store changes to a number of different sounds in working memory (Snyder, Gregg, Weintraub, & Alain, 2012), filtering out sounds that are less relevant for communication and biasing processing towards sounds with high communicative value could be an important step for bootstrapping language learning in the first months of life. As such, detecting small changes in speech sounds may be an important pre-requisite for speech perception and language learning in the real world. Future work could examine whether

a stronger bias for detecting changes in the human voice and speech in complex scenes is related to phonological awareness or word learning by the first year of life.

Beyond their sensitivity to specific sound categories, infants were sensitive to the size of the acoustic change. This was evident when infants looked longer to water sounds changing compared to no change trials, and in their equal looking to changes in water and changes in speech in both experiments. One possibility for the lack of difference between speech and water sounds could relate to scale-invariant or fractal structure that has been found in human speech (Clarke & Voss, 1975) and water sounds (Geffen et al., 2011). However, all of our stimuli were real-world sounds that have natural fluctuations in loudness, pitch, and melodic information, which are the likely drivers of scale-invariant structure in speech, water, and music. Scale-invariance is also evident in natural and environmental sounds, including animal vocalizations (Singh & Theunissen, 2003; Rodriguez, Chen, Read, & Escabi, 2010), which makes it unlikely that scale-invariant structure played a role in looking times for this study.

Infants' sensitivity to the acoustic features of sounds contributed to their looking time differences across all sound types. Looking times were positively correlated with Euclidian distance across participants in Experiment 2 for the speech condition. These results directly align with previous change deafness work suggesting that children and adults use both their knowledge of sound categories and the magnitude of an acoustic change to detect changes across auditory scenes (Gregg & Samuel, 2009; Vanden Bosch der Nederlanden, Hannon, & Snyder, 2016). By 5 months, infants' perceptual abilities to detect changes based on the stimulus features of sounds are already in place. These findings somewhat mirror the development of change detection in complex

visual scenes. Infants looked longer to salient features in child-friendly movies, which declined between 3- and 9-months of age, while looking to faces increased over this period (Frank, Vul & Johnson, 2009). Although infants in the current study already show a bias toward human speech over large acoustic changes, further studies should examine whether infants preferentially attend to large, salient changes early in development across modalities and only begin to attend preferentially to sound categories, like the human voice and speech, with attentional maturity (Frank, Amso, & Johnson, 2014) or with experience.

One limitation of the current study is the small number of simultaneously presented sounds in each scene. As this was the first study to examine change detection in the auditory modality in infancy, we limited scenes to two overlapping sounds. Previous change deafness literature found that children and adults are biased toward the human voice with 4-sound scenes (Vanden Bosch der Nederlanden, Snyder, & Hannon, 2016). Other change deafness studies have examined change detection in scene sizes up to 8 simultaneously presented sounds (Eramudugolla et al., 2005; Gregg, Irsik, & Snyder, 2017). As the acoustic world of young language learners likely contains a large number of sounds in a typical scene, future studies will need to characterize how many sounds infants can process in a scene while still detecting changes between scenes, and whether the bias for speech persists with more sounds in a scene.

Another limitation is the artificial nature of the simultaneous presentation of sounds in our scenes. As this was the first study of change detection in infancy, we based our experimental paradigm on previous literature, but it would be important to use

1 naturally recorded scenes of, for example, a kitchen, a restaurant, a daycare. Natural
 2 recordings make it difficult to control factors like the number of auditory objects in a
 3 scene, the duration of the sound in a scene, the onset and offset of sounds at varying
 4 times through the scene, all elements that are likely to increase or decrease the
 5 difficulty of the scene, further obscuring the contribution of semantic or acoustic
 6 distance factors. Further, sound scenes are seldom experienced without visual
 7 information that would help disambiguate scenes. While our study demonstrates
 8 attentional priorities toward real-world sounds in complex scenes presented
 9 simultaneously, it will be important to characterize these attentional priorities in real-
 10 world audio and audiovisual scenes to better understand attentional biases in the real
 11 world.

12 Our findings show that young infants building their language system not only
 13 have robust and early emerging preferences for listening to speech over other sounds
 14 (e.g., Vouloumanos & Werker, 2007; 2004), but they also detect changes in speech
 15 sounds better than changes of a similar magnitude for other sounds—over and above
 16 the size of acoustic changes—when sounds are presented simultaneously in complex
 17 auditory scenes more typical of the real world.

18

References

- Best, C. T. & Jones, C. (1998). Stimulus-alternation preference procedure to test infant speech discrimination. *Infant Behaviour and Development*, 21, 295. doi: **10.1016/S0163-6383(98)91508-9**
- Bregman, A. S. (1990). Auditory scene analysis: The perceptual organization of sound. Cambridge, MA, US: The MIT Press.
- Cooper, R. P. & Aslin, R. N. (1994). Developmental differences in infant attention to the spectral properties of infant-directed speech. *Child Development*, 65(6), 1663-1677. doi: **10.1111/j.1467-8624.1994.tb00841.x**
- Cousineau, D. (2005). Confidence intervals in within-subject designs: A simpler solution to Loftus and Masson's method. *Tutorials in Quantitative Methods for Psychology*, 1(1), 42-45. doi: **10.20982/tqmp.01.1.p042**
- Cowan, N., Suomi, K., & Morse, P. A. (1982). Echoic storage in infant perception. *Child Development*, 53(4), 984-990. doi: **10.2307/1129138**
- Eramudugolla, R., Irvine, D. R., McAnally, K. I., Martin, R. L., & Mattingley, J. B. (2005). Directed attention eliminates 'change deafness' in complex auditory scenes. *Current Biology*, 15(12), 1108-1113. doi: **10.1016/j.cub.2005.05.051**
- Farroni T, Johnson MH, Menon E, Zulian L, Faraguna D, Csibra G. (2005). Newborns' preference for face- relevant stimuli: Effects of contrast polarity. *Proceedings of the National Academy of Sciences (USA)*, 102,17245–17250.
- Faul, F., Erdfelder, E., Buchner, A., & Lang, A.-G. (2009). Statistical power analyses using G*Power 3.1: Tests for correlation and regression analyses. *Behavior Research Methods*, 41, 1149-1160.

- 1 Fletcher-Watson, S., Leekam, S. R., Benson, V., Frank, M. C., & Findlay, J. M. (2009).
2 Eye-movements reveal attention to social information in autism spectrum
3 disorder. *Neuropsychologia*, 47, 248-257. doi:
4 **10.1016/j.neuropsychologia.2008.07.016**
- 5 Frank, M. C., Amso, D., & Johnson, S. P. (2014). Visual search and attention to faces
6 during early infancy. *Journal of Experimental Child Psychology*, 118, 13-26. doi:
7 **10.1016/j.jecp.2013.08.012**
- 8 Frank, M. C., Vul, E. & Johnson, S. P. (2009). Development of infants' attention to
9 faces during the first year. *Cognition*, 110(2), 160-170. doi:
10 **10.1016/j.cognition.2008.11.010**
- 11 Gliga, T., Elsabbagh, M., Andravizou, A., & Johnson, M. (2009). Faces attract infants'
12 attention in complex displays. *Infancy*, 14(5), 550-562. doi:
13 **10.1080/15250000903144199**
- 14 Gregg, M. K., Irsik, V. C., & Snyder, J. S. (2017). Effects of capacity limits, memory
15 loss, and sound type in change deafness. *Attention, Perception &*
16 *Psychophysics*, 79(8), 2564-2575. doi: **10.3758/s13414-017-1416-4.**
- 17 Gregg, M. K. & Samuel, A. G. (2008). Change deafness and the organizational
18 properties of sounds. *Journal of Experimental Psychology: Human Perception*
19 *and Performance*, 34(4), 974-991. doi: **10.1037/0096-1523.34.4.974**
- 20 Gregg, M. K. & Samuel, A. G. (2009). The importance of semantics in auditory
21 representations. *Attention, Perception & Psychophysics*, 71(3), 607-619. doi:
22 **10.3758/APP.71.3.607**

- 1 Gygi, B., Kidd, G. R., & Watson, C. S. (2007). Similarity and categorization of
2 environmental sounds. *Perception & Psychophysics*, 6(69), 839-855. doi:
3 **10.3758/BF03193921**
- 4 Jusczyk, P. W., Houston, D. M., & Newsome, M. (1999). The beginnings of word
5 segmentation in English-learning infants. *Cognitive Psychology*, 39(3-4), 159-
6 207. doi: **10.1006/cogp.1999.0716**
- 7 Kuhl, P. K. (2004). Early language acquisition: cracking the speech code. *Nature*
8 *Reviews Neuroscience*, 5(11), 831-843. doi: **10.1038/nrn1533**
- 9 Lim, S. J., Lacerda, F., & Holt, L. L. (2015). Discovering functional units in continuous
10 speech. *Journal of Experimental Psychology: Human Perception and*
11 *Performance*, 41(4), 1139-1152. doi: **10.1037/xhp0000067**
- 12 Mitroff, S. R., Simons, D. J., & Levin, D. T. (2004). Nothing compares 2 views: change
13 blindness can occur despite preserved access to the changed information.
14 *Perception & Psychophysics*, 66(8), 1268-1281. doi: **10.3758/BF03194997**
- 15 Newman, R. S. & Jusczyk, P. W. (1996). The cocktail party effect in
16 infants. *Perception & Psychophysics*, 58(8), 1145-1156. doi:
17 **10.3758/BF03207548**
- 18 Newman, R., Ratner, N. B., Jusczyk, A. M., Jusczyk, P. W., Dow, K. A. (2006). Infants'
19 early ability to segment the conversation speech signal predicts later language
20 development. A retrospective analysis. *Developmental Psychology*, 42(4), 643-
21 655. doi: **10.1037/0012-1649.42.4.643**
- 22 Oakes, L.M., Sperka, D. J., & Cantrell, L. (2015). Habit 2. *Unpublished software*.
23 *Center for Mind and Brain, University of California, Davis.*

- 1 Rensink, R. A., O'Regan, J. K., & Clark, J. J. (1997). To see or not to see: The need
2 for attention to perceive changes in scenes. *Psychological Science*, 8(5), 368-
3 373. doi: **10.1111/j.1467-9280.1997.tb00427.x**
- 4 Ro, T., Russell, C., & Lavie, N. (2001). Changing faces: A detection advantage in the
5 flicker paradigm. *Psychological Science*, 12(1), 94-99. doi: **10.1111/1467-**
6 **9280.00317**
- 7 Saffran, J. R., Aslin, R. N., & Newport, E. L. (1996). Statistical learning by 8-month-old
8 infants. *Science*, 274(5294): 1926-1928. doi: **10.1126/science.274.5294.1926**
- 9 Scerif, G. (2010). Attention trajectories, mechanisms and outcomes: at the interface
10 between developing cognition and environment. *Developmental Science*, 13(6),
11 805-812. doi: **10.1111/j.1467-7687.2010.01013.x**
- 12 Shultz, S. & Vouloumanos, A. (2010). Three-month-olds prefer speech to other
13 naturally occurring signals. *Language Learning and Development*, 6, 241-257.
14 doi: **10.1080/15475440903507830**
- 15 Simion F, Macchi Cassia V, Turati C, Valenza E. (2001). The origins of face
16 perception: specific versus non- specific mechanisms. *Infant and Child*
17 *Development*,10,59–65.
- 18 Simons, D. J. (2000). Current approaches to change blindness. *Visual Cognition*, 7(1-
19 3), 1-15. doi: **10.1080/135062800394658**
- 20 Simons, D. J. & Rensink, R. A. (2005). Change blindness: past, present, and future.
21 *Trends in Cognitive Sciences*, 9(1), 16-20. doi: **10.1016/j.tics.2004.11.006**

- 1 Snyder, J.S. & Gregg, M.K. (2011). Memory for sound, with an ear toward hearing in
2 complex auditory scenes. *Attention, Perception, and Psychophysics*, 73(7), 1993-
3 2007.
- 4 Snyder, J. S., Gregg, M. K., Weintraub, D. M., & Alain, C. (2012). Attention,
5 awareness, and the perception of auditory scenes. *Frontiers in Psychology*,
6 3, Article ID 15. doi: **10.3389/fpsyg.2012.00015**
- 7 Swingle, D. (2009). Onsets and codas in 1.5-year-olds' word recognition. *Journal of*
8 *Memory and Language*, 60(2), 252-269. doi: **10.1016/j.jml.2008.11.003**
- 9 Thiessen, E. D., Hill, E. A., & Saffran, J. R. (2005). Infant-directed speech facilitates
10 word segmentation. *Infancy*, 7(1), 53-71. doi: **10.1207/s15327078in0701_5**
- 11 Vanden Bosch der Nederlanden, C. M., Snyder, J. S., & Hannon, E. E. (2015).
12 Everyday musical experience is sufficient to perceive the speech-to-song illusion.
13 *Journal of Experimental Psychology: General*, 144(2), e43-e49. doi:
14 **10.1037/xge0000056**
- 15 Vanden Bosch der Nederlanden, C. M., Snyder, J. S., & Hannon, E. E. (2016).
16 Children use object-level category knowledge to detect changes in complex
17 auditory scenes. *Developmental Psychology*, 52(11), 1867-1877. doi:
18 **10.1037/dev0000211**
- 19 Vanden Bosch der Nederlanden, C. M., Zaragoza, C., Rubio-Garcia, A., Clarkson,
20 E., & Snyder, J. S. (2018). Change detection in complex auditory scenes is
21 predicted by auditory memory, pitch perception, and years of musical training.
22 *Psychological Research*, 1-17. doi: **10.1007/s00426-018-1072-x**

- 1 Vitevitch, M (2003). Change deafness: The inability to detect changes between two
2 voices. *Journal of Experimental Psychology: Human Perception and*
3 *Performance*, 29(2), 333-342.
- 4 Vouloumanos, A. & Werker, J. F. (2004). Tuned to the signal: the privileged status of
5 speech for young infants. *Developmental Science*, 7(3), 270-276. doi:
6 **10.1111/j.1467-7687.2004.00345.x**
- 7 Vouloumanos, A. & Werker, J. F. (2007). Listening to language at birth: evidence for a
8 bias for speech in neonates. *Developmental Science*, 10(2), 159-164. doi:
9 **10.1111/j.1467-7687.2007.00549.x**
- 10 Vouloumanos, A. & Waxman, S. R. (2014). Listen up! Speech is for thinking during
11 infancy. *Trends in Cognitive Neuroscience*, 18(12), 642-646. doi:
12 **10.1016/j.tics.2014.10.001**
- 13 Waxman, S. R. & Markow, D. B. (1995). Words as invitations to form categories:
14 Evidence from 12- to 13-month-old infants. *Cognitive Psychology*, 29(3), 257-
15 302. doi: **10.1006/cogp.1995.1016**
- 16 Werker, J. F., Cohen, L. B., Lloyd, V. L., Casasola, M., & Stager, C. L. (1998).
17 Acquisition of word-object associations by 14-month-old infants. *Developmental*
18 *Psychology*, 34(6), 1289-1309. doi: **10.1037/0012-1649.34.6.1289**
- 19 Werker, J. F. & Tees, R. C. (1984). Cross-language speech perception: Evidence for
20 perceptual reorganization during the first year of life. *Infant Behavior &*
21 *Development*, 7(1), 49-63. doi: **10.1016/S0163-6383(84)80022-3**

- 1 Werker, J. F., Yeung, H. H., & Yoshida, K. A. (2012). How do infants become experts
- 2 at native-speech perception? *Current Directions in Psychological Science*, 21(4),
- 3 221-226. doi: **10.1177/0963721412449459**
- 4

1

Appendix

Set 1				Set 2			
Trial	Scene 1	Scene 2	Sound Change Category	Trial	Scene 1	Scene 2	Sound Change Category
1	Female Speaking Bubbling Water	Female Speaking Bubbling Water	No Change	10	Male Speaking Draining Water	Male Speaking Draining Water	No Change
2	Female Speaking Bubbling Water	Cello Bubbling Water	Speech Change	11	Male Speaking Draining Water	Tuba Draining Water	Speech Change
3	Female Speaking Bubbling Water	Female Speaking Cello	Water Change	12	Male Speaking Draining Water	Male Speaking Tuba	Water Change
4	Female Speaking Cello	Female Speaking Cello	No Change	13	Male Speaking Tuba	Male Speaking Tuba	No Change
5	Female Speaking Cello	Bubbling Water Cello	Speech Change	14	Male Speaking Tuba	Draining Water Tuba	Speech Change
6	Female Speaking Cello	Female Speaking Bubbling Water	Music Change	15	Male Speaking Tuba	Male Speaking Draining Water	Music Change
7	Bubbling Water Cello	Bubbling Water Cello	Water Change	16	Draining Water Tuba	Draining Water Tuba	Water Change
8	Bubbling Water Cello	Female Speaking Cello	Water Change	17	Draining Water Tuba	Male Speaking Tuba	Water Change
9	Bubbling Water Cello	Bubbling Water Female Speaking	Music Change	18	Draining Water Tuba	Draining Water Male Speaking	Music Change

2