

Invariance in pitch perception

Malinda J. McPherson^{1,2,3}, Josh H. McDermott^{1,2,3,4}

1 Department of Brain and Cognitive Sciences, MIT

2 Program in Speech and Hearing Biosciences and Technology, Harvard University

3 McGovern Institute for Brain Research, MIT

4 Center for Brains Minds and Machines, MIT

ABSTRACT

Information in speech and music is often conveyed through changes in fundamental frequency (f_0), the perceptual correlate of which is known as “pitch”. One challenge of extracting this information is that such sounds can also vary in their spectral content due to the filtering imposed by a vocal tract or instrument body. Pitch is envisioned as invariant to spectral shape, potentially providing a solution to this challenge, but the extent and nature of this invariance remain poorly understood. We examined the extent to which human pitch judgments are invariant to spectral differences between natural sounds. Listeners performed up/down and interval discrimination tasks with spoken vowels, instrument notes, or synthetic tones, synthesized to be either harmonic or inharmonic (lacking a well-defined f_0). Listeners were worse at discriminating pitch across different vowel and instrument sounds compared to when vowels/instruments were the same, being biased by differences in the spectral centroids of the sounds being compared. However, there was no interaction between this effect and that of inharmonicity. In addition, this bias decreased when sounds were separated by short delays. This finding suggests that the representation of a sound’s pitch is itself unbiased, but that pitch comparisons between sounds are influenced by changes in timbre, the effect of which weakens over time. Pitch representations thus appears to be relatively invariant to spectral shape. But relative pitch judgments are not, even when spectral shape variation is naturalistic, and when such judgments are based on representations of the f_0 .

1 INTRODUCTION

A central challenge for our perceptual systems is that we must often make judgments about one variable amid variation across other variables (Carruthers et al., 2015; DiCarlo & Cox, 2007; Liu et al., 2019; Sharpee et al., 2011). For example, object shape must be estimated across changes in lighting and pose, and spoken words recognized across variation in the voice of the speaker, manner of speaking, and listening environment.

Another instance of this challenge can be found in pitch perception. ‘Pitch’ is typically defined within auditory science as the perceptual correlate of a sound’s fundamental frequency (f_0) (Plack et al., 2005). Natural sounds are often harmonic, such that their frequencies are integer multiples of an f_0 . But the relative amplitudes of these harmonics can vary, such that sounds with the same f_0 can have different spectra. This variation is intrinsic to the source-filter generative model of speech and instrument sounds (Fletcher & Rossing, 2010; Stevens, 2000). The source is characterized in part by its fundamental frequency (f_0). The filter is mediated by the vocal tract or instrument body, which have resonances that amplify some frequencies and attenuate others.

Pitch is believed to be somewhat invariant to the spectral shape imposed by such filters (de Cheveigne, 2010; Semal & Demany, 1991). One example of this invariance is the ‘missing fundamental illusion’. When the lowest frequency component is removed from a harmonic tone, the f_0 of the tone remains the same, and listeners are thought to be able to hear this invariant property of the tones even though the spectrum has changed (Licklider, 1954). Many studies have explored variants of this phenomenon with synthetic tones, testing pitch matching or discrimination of tones that vary in their spectra. Human pitch discrimination exhibits some robustness to synthetic spectral variation, in that listeners are above chance at discriminating tones that have non-overlapping sets of harmonics (Micheyl & Oxenham, 2004; Moore et al., 1992; Singh & Hirsh, 1992). However, discrimination is nonetheless impaired compared to when tones have the same set of harmonics (Allen & Oxenham, 2014; Melara & Marks, 1990; Micheyl & Oxenham, 2004; Moore et al., 1992; Moore & Glasberg, 1990; Russo & Thompson, 2005; Singh & Hirsh, 1992; Warrier & Zatorre, 2002). A priori it seemed plausible that these impairments could result from spectral differences that are relatively extreme compared to what occurs in natural sounds.

The effect of naturally occurring spectral variation (as is present between different syllables or instruments) on pitch perception has remained relatively unexplored. Two previous studies compared pitch judgments of pairs of notes played on same vs. different instruments (Vurma et al., 2011; Zarate et al., 2013), finding slight deficits in performance when notes were from different instruments. But these studies left open the basis of the observed deficit, as well as the residual invariance. And to our knowledge, there are no existing measurements of pitch judgments for speech sounds with varied timbre.

One might suppose that the ability to discriminate changes in f_0 despite concurrent changes in spectral shape would be mediated by representations of the f_0 that are invariant to spectral shape. However, when the f_0 of a sound changes, there are corresponding changes in the frequencies of individual harmonics, which shift up or down with the f_0 . In many conditions listeners appear to use the harmonics rather than the f_0 to judge whether one sound is higher or lower than another (Faulkner, 1985; McPherson & McDermott, 2018; McPherson & McDermott, 2020; Micheyl et al., 2010; Moore & Glasberg, 1990), and it was thus unclear whether any spectral invariance of pitch judgments would depend on representations of the f_0 . One tool to test for representations of the f_0 is to compare performance with harmonic sounds to that with inharmonic sounds (which lack a single f_0 in the range of audible pitch). A role for f_0 representations in spectral invariance has been suggested by discrimination advantages for harmonic tones over inharmonic tones when the tones being compared contain distinct sets of harmonics (Micheyl et al., 2010; Moore & Glasberg, 1990). However, these studies left it unclear whether this advantage

would be observed for naturally occurring spectral variation, as when comparing the pitch of different vowels in a speech utterance, or the pitch of different instruments.

The goals of this paper were to assess 1) the extent to which human discrimination is invariant to naturally occurring spectral differences, 2) whether any such invariance is dependent on representations of f_0 and 3) to what extent any limits on invariance reflect imperfect invariance of pitch representations themselves vs. biases in the judgments that operate on those representations. To examine the first question, we measured f_0 discrimination with the same or different spoken vowels, as well as with notes played on the same or different musical instruments. If pitch perception is invariant to real-world spectral changes, pitch discrimination performance should be similar for same vs. different musical instruments and vowels. To answer the second question, we compared performance for harmonic and inharmonic stimuli. Our assumption was that if a task relies on representations of the f_0 , performance should be impaired with inharmonic stimuli (Faulkner, 1985; McPherson & McDermott, 2018; McPherson & McDermott, 2020; Micheyl et al., 2010; Moore & Glasberg, 1990). Specifically, if spectral invariance is mediated by f_0 -based pitch, we predicted that there would be significant harmonic advantages when discriminating sounds that differ in their spectra. To answer the third question, we measured pitch discrimination across short time delays, leveraging the apparent robustness of f_0 representations over time (McPherson & McDermott, 2020). Our hypothesis was that if limits to spectral invariance arise only at a decision stage, effects of the spectrum on discrimination might decrease across a delay, as the representation of the spectrum might degrade more rapidly than that of the f_0 .

We found that discrimination was only partially invariant to real-world differences in spectral shape between different vowels or instruments. Judgments were biased depending on whether the spectrum shifted congruently or incongruently with the f_0 , consistent with previous studies with synthetic tones (Micheyl & Oxenham, 2004; Moore & Glasberg, 1990; Singh & Hirsh, 1992). However, invariance was similarly limited for harmonic and inharmonic tones, indicating that invariance does not depend on representations of f_0 . Invariance increased when we introduced delays between tones, suggesting that representations of pitch and spectral shape decay at different rates, and are thus independent. The observed biases appear to occur at a comparison stage that cannot separate the effects of pitch and timbre. The results suggest that pitch representations are relatively invariant to spectral shape. But relative pitch judgments are not, even when the variation in spectral shape is naturally occurring, and irrespective of whether such judgments are based on representations of the f_0 .

2 METHODS AND MATERIALS

This section contains aspects of the methods that were shared by two or more experiments.

2.1 Ethics

All experiments were approved by the Committee on the use of Humans as Experimental Subjects at the Massachusetts Institute of Technology and were conducted with the informed consent of the participants.

2.2 Audio presentation and procedure for online experiments (Experiments 1-2, 4-5)

Experiments 1, 2, 4, and 5 were completed online due to the COVID-19 pandemic. Online psychoacoustic experiments sacrifice control over absolute sound presentation levels and spectra, but we have repeatedly found that online experiments replicate results obtained in controlled laboratory conditions provided that modest steps are taken to help ensure reasonable sound quality and compliance with task instructions (Kell et al., 2018; McPherson et al., 2020; McPherson et al., 2021; McPherson & McDermott, 2020; McWalter & McDermott, 2019; Traer et al., 2021; Woods & McDermott, 2018). Experiments were conducted on the Amazon Mechanical Turk platform. We limited participation to individuals with US-based IP addresses. Before beginning the experiment, potential participants were consented and instructed to wear headphones and ensure they were in a quiet location. They then used a calibration sound (1.5 seconds of speech-shaped noise) to set their audio presentation volume to a comfortable level. The experimental stimuli were normalized to 6 dB below the level of the calibration sound to ensure that they were likely to be audible without being uncomfortably loud. Participants then completed a brief screening test designed to help ensure they were wearing earphones or headphones (Woods et al., 2017). If they passed this screening, they could continue to the main experiment. To incentivize good performance, online participants received a compensation bonus proportional to the number of correct trials. All online experiments began with a set of screening questions that included a question asking the participant if they had any hearing loss. Anyone who indicated any known hearing loss was excluded from the study. Across all the online experiments in this paper, 3.5% of the 1052 participants who initially enrolled self-reported hearing loss and were excluded (all but four of these individuals also failed the headphone screening, and would have been excluded regardless). All participants whose data were analyzed thus self-reported normal hearing.

For technical reasons all stimuli for online experiments were generated ahead of time and were stored as .wav files on a university server, from which they could be loaded during the experiments. 20 sets of stimuli were pre-generated for each experiment, and participants only heard stimuli from one of these sets, randomly assigned.

2.3 Inharmonic stimuli

In order to make synthetic tones, spoken syllables, or musical instrument notes inharmonic, the frequency of each harmonic, excluding the fundamental, was perturbed (jittered) by an amount chosen randomly from a uniform distribution, $U(-.5, .5)$. These jitter values were multiplied by the f_0 of the tone and added to the frequency of the respective harmonic. For example, if the f_0 was 200 Hz and a jitter value of -0.39 was selected for the second harmonic; its frequency would be set to 322 Hz ($200 \cdot 2 + 200 \cdot -.39$, Figure 1a). To minimize salient differences in beating, jitter values were constrained via rejection sampling such that adjacent harmonics were always separated by at least 30 Hz. Jitter values were generated for each harmonic in succession, beginning with the second, subject to the 30 Hz constraint. A different jitter pattern was chosen randomly for each trial, but the same jitter pattern was applied to the two tones/notes/vowels within a trial.

2.4 Distortion products

Distortion products can in principle explain differences in performance for harmonic and inharmonic stimuli (because they should be stronger in harmonic stimuli). However, in all cases in which we

compared task performance for harmonic and inharmonic stimuli (Experiments 1, 2, and 5), the stimuli contained all lower harmonics, including the fundamental. Because distortion products are typically substantially lower in level than the stimulus components that generate them (Norman-Haignere & McDermott, 2016; Pressnitzer & Patterson, 2001), they are unlikely to be detectable in stimuli that contain all lower harmonics, and thus are unlikely to account for any differences between performance for harmonic and inharmonic stimuli. We thus did not include masking noise in the stimuli.

2.5 Feedback

Feedback was given for all tasks except Experiments 3 (where we anticipated chance or below-chance performance for most of the Inharmonic conditions, and we did not want participants to get discouraged).

2.6 Statistics

Data distributions were evaluated for normality by visual inspection and parametric statistics were used across all conditions. Unless otherwise noted we tested hypotheses using repeated measures ANOVAs. All analysis was completed in MATLAB (version 2020b). Bayesian statistics were used to evaluate null results (*JASP*, 2020).

3 RESULTS

3.1 Experiment 1: Pitch discrimination with same vs. different vowels

3.1.1 Purpose and procedure

We began by testing up/down discrimination of same vs. different vowels. During each trial, participants heard two vowels and judged whether the second was higher or lower in pitch than the first, indicating their response by clicking one of two buttons ('higher' or 'lower', Figure 1b). The conditions resulted from fully crossing three variables: Same vs. Different vowels, Harmonic vs. Inharmonic, and 4 frequency/f₀ differences (.33, 1, 3, and 9 semitones). Conditions were randomly intermixed during the experiment, and participants completed 30 trials per condition.

Our first hypothesis was that human listeners are invariant to natural spectral differences, in which case there should be no significant difference between discrimination for same vs. different vowels. Our second hypothesis was that any spectral invariance (ability to discriminate pitch despite differences in the spectra) would be mediated by representations of the f₀, which would be indicated by a harmonic advantage specific to (or greater when) discriminating different vowels.

3.1.2 Stimuli

During each trial, participants were asked to compare two vowels, each spoken by the same speaker to maximize ecological validity (we envisioned the task as potentially tapping into the mechanisms underlying the perception of the prosodic contour of an utterance). We aimed to choose a set of pairings for which the two vowels in each pair varied substantially in their spectral shape (Figure 1c-d; example same vs. different vowel spectra).

To find vowel pairs that differed in their spectral shape, we compared the excitation patterns across pairs of different vowels with the same f₀. Excitation patterns were calculated by passing waveforms (each RMS normalized to the same level) through a gammatone filter bank (Slaney, 1998) approximating the frequency selectivity of the cochlea. The filter bank was implemented with Malcolm Slaney's Auditory Toolbox, with the lowest center frequency set to 50 Hz, the number of channels set to 64, and the sampling rate set to 16 kHz. The excitation pattern was calculated from a magnitude spectrogram generated from these filter outputs (the RMS amplitude within .025 s bins; hop size of .01 s) by averaging the spectrogram over time. The spectral difference between two vowels was calculated as the sum over frequency of the absolute value of the dB differences between the two excitation patterns (Figure 1c).

Vowels were selected from the Hillenbrand vowel set (Hillenbrand et al., 1995). We included the vowels /I, I, æ, a, ɔ, ε, u, ʌ, ʊ, ɜ~/ in our analysis, omitting the two diphthongs in the Hillenbrand set (/ou/ and /eɪ/) because their spectra are not fixed across their duration. We calculated the spectral difference between every pair of the vowels (45 pairs) for every speaker (45 male speakers and 48 female speakers). Before comparing the spectra of each vowel in a pair, we pitch-shifted the two vowels to match their mean f₀ (we pitch-shifted each of the two vowels by an equal and opposite amount, using STRAIGHT, described below). We settled on the set of 10 vowel pairings that maximized the average spectral difference between vowels of the same speaker. We did not otherwise constrain the vowel pairings, so some vowel tokens are over-represented (see Supplementary Table 1 for the final set of vowels).

For every 'Different Vowel' condition, we used each of the 10 vowel pairings 3 times, randomly sampling from the speakers in the set. 15 of these 30 vowel sets were selected from female speakers, and 15 from male speakers. The presentation order of the two vowels within a trial was randomized.

For 'Same Vowel' conditions, we chose a single vowel recording to use for both intervals of the trial. We sampled a set of 30 vowels that included all 20 of the vowels from the 10-vowel pairings set along with

an additional 10 vowels that were sampled randomly from the set of 20 without replacement. We again sampled randomly from the speakers in the set, balancing for speaker gender.

To make the final stimuli, vowels were pitch-shifted and resynthesized to be either harmonic or inharmonic using the STRAIGHT analysis and synthesis method (Kawahara, 2006; Kawahara & Morise, 2011; McDermott et al., 2012). STRAIGHT decomposes a recording of speech or instruments into periodic and aperiodic excitation and a time-varying filter. If the periodic excitation is modeled sinusoidally, one can alter the frequencies of individual harmonics, and then recombine them with the unaltered aperiodic excitation and filter to generate harmonic or inharmonic speech. This manipulation leaves the spectral shape of the speech largely intact, and a previous study suggests that intelligibility of inharmonic speech in quiet is comparable to that of harmonic speech (Popham et al., 2018).

We pitch-shifted vowels to achieve the desired f_0 difference, but did not otherwise modify the pitch contours of the vowels. To perform the pitch adjustments, we first calculated the combined mean f_0 across both vowels. We then shifted both vowels equal distances towards or away from their combined mean f_0 such that the mean pitch difference between the two vowels matched the exact pitch difference needed for the respective condition. This ensured that the pitch of both vowels was shifted as little as possible from the original pitch, and by equal amounts. We tested four different step sizes: .33, 1, 3 and 9 semitones.

Because we did not pitch-flatten or otherwise modify the pitch contours of the vowels, it seemed important to control for the fact that in the Different Vowel case, the two selected vowels would have different pitch fluctuations. We addressed this possible confound by giving the two intervals in a Same Vowel trial different pitch contours. This was achieved by replacing the pitch contour of one of the two vowels in the 'Same Vowel' condition with the pitch contour from what would have been the other vowel, were it a Different Vowel condition. For example, if the vowel in a Same Vowel trial was the vowel /æ/ drawn from the pair /æ/ and /u/, we would use a single recording of the /æ/ vowel from one speaker for both intervals in the trial, but would replace the f_0 contour for one of the intervals with that from the speaker's /u/ vowel. In this way, the pitch variability within trials was matched between Same Vowel and Different Vowel conditions.

The resynthesized vowels were truncated at 400 ms, windowed with a 15 ms half-Hanning window, and presented in sequence (with no time delay between vowels). Sounds were sampled at 16 kHz (the sampling rate of the Hillenbrand vowel recordings).

3.1.3 Participants

127 participants passed the headphone check and enrolled in Experiment 1 online. 45 had overall performance (averaged across all conditions) below 55% correct, and were removed from further analysis, leaving 82 participants whose data are reported here. 26 identified as female, 56 as male, 0 as non-binary. The average age of these participants was 38.0 years (S.D.=10.4). 34 participants had four or more years of musical training (self-reported), with an average of 11.4 years (S.D.=8.4).

We based our initial target sample size on the same pilot experiment and power analysis used for Experiment 2 (described below), which showed that we needed 7 participants to have a 95% chance of observing an effect of same vs. different instruments at a significance level of .01. As we observed a null effect of harmonicity for small f_0 differences after collecting our initial data, we continued data collection until Bayesian statistics converged on support for or against the null hypothesis.

3.1.4 Results and discussion

Our first hypothesis was that participants are invariant to spectral differences between vowels, such that we would observe similar discrimination across same and different vowels. Instead, we found that discrimination was worse when participants compared different vowels than when they compared instances of the same vowel (main effect of Same vs. Different vowels, $F(1,81)=78.76$, $p<.0001$, $\eta_p^2=.49$, Figure 1e). However, participants were nonetheless well above chance when discriminating the pitch of different vowels. Listeners thus exhibited some degree of spectral invariance.

Our second hypothesis was that harmonicity would be critical to any observed spectral invariance, such that there might be a harmonic advantage in the different vowel condition but not the same vowel condition. Contrary to this hypothesis, there was no interaction between the effects of Same vs. Different vowels and Harmonicity ($F(1,81)=0.20$, $p=.66$, $\eta_p^2=.002$, the $BF_{incl}=0.76$, specifying a multivariate Cauchy prior on the effects, provided anecdotal support for the null hypothesis (*JASP*, 2020)). In particular, for the three smallest step sizes, there was no difference between Harmonic and Inharmonic performance for either the Same Vowel or Different Vowel conditions (no main effect in either case; Same Vowel: $F(1,81)=0.37$, $p=.55$, $\eta_p^2=.005$, Different Vowel: $F(1,81)=0.66$, $p=.42$, $\eta_p^2=.008$). The Bayes factors in each case (BF_{incl} of .13 and .11 for Same and Different Vowels, respectively) provided strong support for the null hypotheses (*JASP*, 2020).

However, there were clear differences between performance with Harmonic and Inharmonic sounds for both vowel conditions when the step size was large. At a 9 semitone difference there was a highly significant effect of Harmonicity ($F(1,81)=104.91$, $p<.0001$, $\eta_p^2=.56$; this drove an overall effect of Harmonicity across all step sizes (i.e. the f_0 difference, or nominal f_0 difference in the inharmonic case, between the two sounds on a trial, in semitones), $F(1,81)=18.94$, $p<.0001$, $\eta_p^2=.19$). This effect with large intervals resulted in a highly significant interaction between step size and Harmonicity ($F(3,243)=50.71$, $p<.0001$, $\eta_p^2=.38$). This result is similar to one that we previously found for synthetic tones with a fixed filter (replotted for convenience in Supplementary Figure 1), and suggests the effect has relevance for real-world listening conditions.

All together, these results suggest that participants are not fully invariant to changes in the spectrum, and that the invariance they have is not mediated by a representation of f_0 . However, representations of f_0 appear to be helpful in discriminating large pitch differences. For large pitch differences it may be difficult for listeners to match the frequency components of one inharmonic tone to those of the next inharmonic tone, as is necessary to determine the step from a representation of the fine spectrum (Figure 1f). Having access to the f_0 (as is available in harmonic tones) may help listeners resolve these ambiguities.

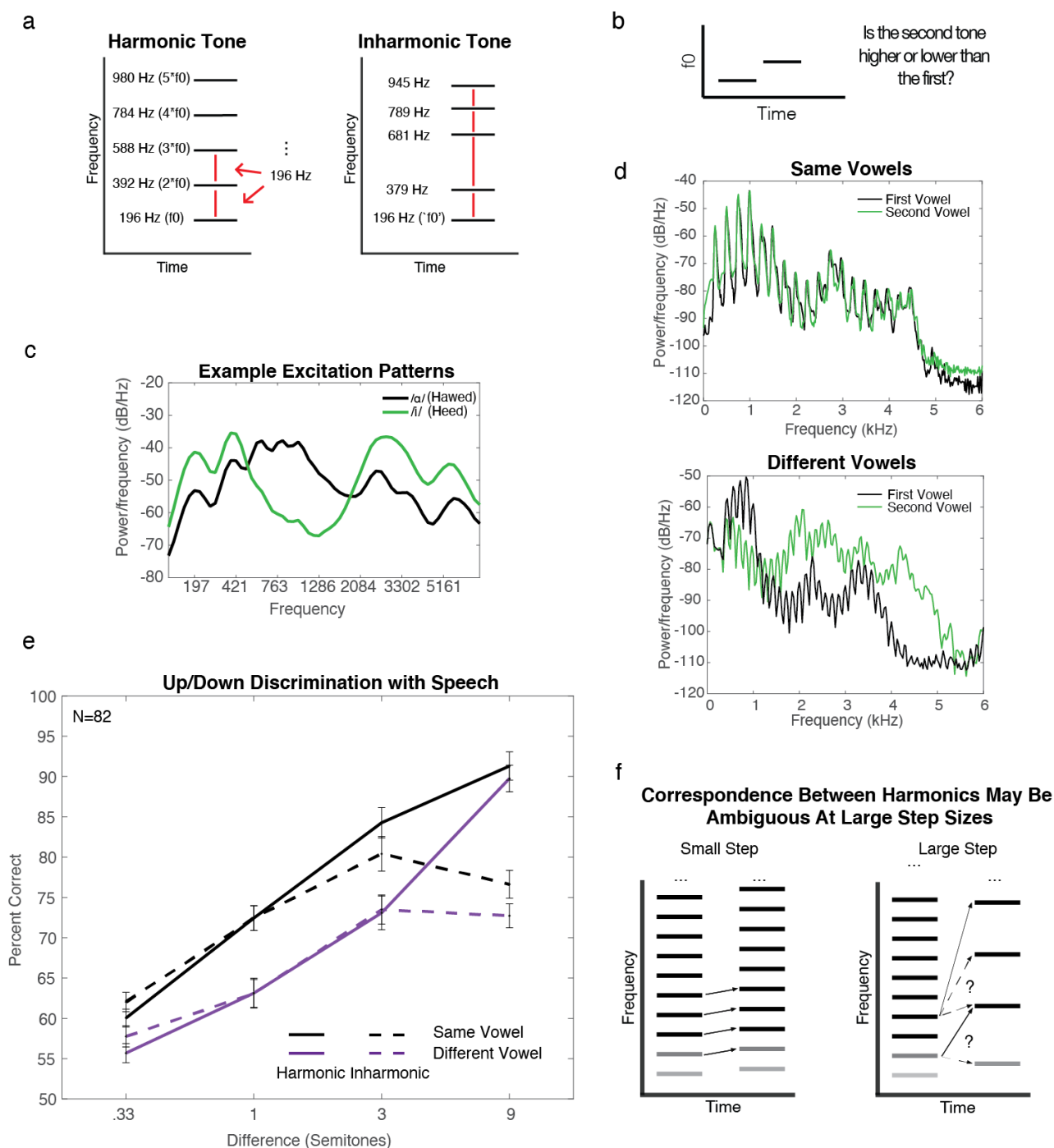


Figure 1: Stimuli and Results for Experiment 1, Pitch discrimination with same vs. different vowels

a. Schematic spectrogram of a 196 Hz Harmonic tone (G3 in the Western classical scale), and inharmonic tone with a nominal f_0 of 196 Hz. The inharmonic tone was generated by altering the frequencies of all harmonics above the f_0 , by sampling a jitter value from the distribution $U(-0.5, 0.5)$, multiplying that jitter by the f_0 , then adding the resulting value to the frequency of the respective harmonic, constraining adjacent harmonics to be separated by at least 30 Hz (via rejection sampling) in order to avoid salient beating. b. Task for all experiments. Participants heard two sounds and judged which was higher in pitch. c. Example excitation patterns for two vowels (x-axis scaled in Equivalent Rectangular Bandwidths based on Glasberg and Moore (1990)). Excitation patterns contain peaks for individual harmonics, but also reveal the

differences in spectral shape between the vowels d. Example stimulus power spectra from two exemplars of the same vowel (top), and two different vowels (bottom). e. Results from Experiment 1. Error bars denote standard error of the mean. f. Schematic illustrating potential difficulty in matching harmonics between tones with a fixed filter that are separated by a large step size. Gray level of line segments denotes amplitude (lower harmonics are attenuated by the filter). To facilitate inspection of the ambiguity in matching, only lower portion of spectrum is shown. When two tones have similar f_0 s, the correspondence between harmonics is relatively unambiguous (left), but when the f_0 change is large enough the correspondence between harmonics is ambiguous (right), impairing performance in conditions in which listeners base judgments on shifts in harmonics rather than the f_0 .

3.2 Experiment 2: Pitch discrimination with notes played on the same vs. different instruments

3.2.1 Purpose and procedure

The procedure and hypotheses for Experiment 2 were analogous to those of Experiment 1. However, instead of hearing pairs of vowels, participants heard pairs of instrument notes that were resynthesized to be either harmonic or inharmonic. Participants again completed 30 trials per condition.

3.2.2 Stimuli

Each trial contained two notes resynthesized from recordings of instruments. To approximately replicate the maximal spectral difference between instruments that a listener might hear in Western music, we analyzed the spectra of different instruments to identify a subset that differed the most in their spectra, using the same excitation pattern analysis procedure used with vowels in Experiment 1. We then selected pairs of instruments from this subset for the Different Instrument trials, and single instruments from this subset for Same Instrument trials.

Recordings were drawn from the RWC Music Database of Musical Instrument Sounds (Goto et al., 2003). We excluded instruments whose recordings did not include all notes in the Western Scale within the approximately 200-400 Hz range we were planning to test in the experiment. This criterion excluded several non-Western instruments (Sho, Koto, etc.) and instruments with low or high pitch ranges (contrabass, soprano recorders etc.). Percussion instruments were also excluded because they are naturally inharmonic. These exclusions left us with 25 instruments (see Supplementary Table 2). We then searched through all instrument pairings to find a set of instruments that maximized variability in timbre.

For each pairing, we compared the excitation patterns of notes with the same f_0 , for all notes between G3 and G4 (196-392 Hz). Sets of 5 instruments were then ranked by the pairwise spectral difference summed across all pairs of instruments within the set, and the combination that maximized this quantity was selected. We settled on an instrument set size of 5 because increasing the set size above 5 decreased the average variability between instruments. Because plucked instruments are not perfectly harmonic, it seemed wise not to have too many of them (even though all sounds were resynthesized to be harmonic for the experiment), so we constrained the sets of instruments so that they could only contain one plucked instrument. The final set included baritone saxophone, cello, oboe, pipe organ, and ukulele (Figure 2a). There were 10 possible pairs of instruments that could be drawn from this set. We adopted this procedure (in contrast to that used in Experiment 1, where we selected 10 pairs of vowels with substantial spectral differences) because the spectral differences between instruments were sufficiently pronounced that it was possible to find a set of 5 for which all pairings were between instruments with distinct timbres.

To generate individual trials, the first note was randomly selected from a uniform distribution over the notes in a Western classical chromatic scale between G3 and G4 (196-392 Hz). A recording of this note, from a randomly selected instrument, was chosen as the source for the first note in the trial. The second

note was drawn from either the same instrument or a different instrument, depending on the condition. If different, the second instrument was drawn at random from the remaining 4 instruments. The recording whose f_0 was closest to that needed for the intended frequency/ f_0 difference was chosen to generate the second stimulus. For differences less than one semitone, a note one semitone above or below the first note was chosen as the source for the second note in the trial. We tested five different step sizes: .11, .33, 1, 3 and 9 semitones. We included one step that was smaller than those in Experiment 1 because we pitch flattened the instrument stimuli, making the task easier.

The two notes were then modified using the STRAIGHT analysis and synthesis method (Kawahara, 2006; Kawahara & Morise, 2011; McDermott et al., 2012), in the same way we modified vowels in Experiment 1, to conform to the stimulus parameters for the condition. We confirmed by eye that the spectral envelopes estimated by STRAIGHT for the instruments in the set were reasonable in all cases. The same spectral envelopes were used for harmonic and inharmonic notes, such that any deviations from the true spectral envelopes of the instruments would not influence the results. After analysis, the f_0 contours were pitch-flattened to remove any vibrato and shifted to ensure that the f_0 differences were exact. The notes were then resynthesized with harmonic or inharmonic excitation (Figure 2b). We shifted both notes by the same amount to achieve the desired f_0 differences (in other words, when the f_0 difference of the source notes was less than the f_0 difference for the condition, the lower note's f_0 was decreased and the higher note's f_0 was increased, and vice versa). As a result of the resynthesis, the small deviations from perfect harmonicity found in natural instruments were removed for the harmonic conditions (Fletcher, 1999). The resynthesized notes were truncated at 400 ms, windowed with a 15 ms half-Hanning window, and presented sequentially with no delay between notes. Sounds were sampled at 16 kHz.

3.2.3 Participants

105 participants passed our headphone check and completed Experiment 2 online. 19 had overall performance (averaged across all conditions) below 55% correct and were removed from further analysis. Of the 86 remaining participants, 25 identified as female, 61 as male, 0 as non-binary. The average age of these participants was 37.4 years (S.D.=10.7). 32 participants had four or more years of musical training (self-reported), with an average of 10.6 years (S.D.=10.0).

The sample size was ultimately determined by the desire to provide evidence for or against the null hypothesis that discrimination was the same for harmonic and inharmonic stimuli. We initially performed a power analysis with pilot data. The pilot experiment had slightly different stimuli (same sets of instruments, but instrument notes were high-passed filtered and had low-pass noise added to mask the f_0 component), did not include a 9-semitone step size, and was run in the lab. This pilot experiment showed a main effect of same vs. different instruments ($\eta_p^2=.64$), but no effect of harmonicity. We thus wanted to be well-powered to test for an effect of same vs. different instrument, but also to provide evidence for a null effect of harmonicity should it hold. The power analysis (using G*Power) indicated that we would need only 7 participants to observe the effect of same vs. different instruments at a significance level of .01, 95% of the time (Faul et al., 2007). We aimed to have at least this many musicians and non-musicians, to be able to analyze the groups separately (see *Effects of Musicianship* section below). After replicating the null effect of harmonicity for small f_0 differences in the first 7 participants who completed the experiment, we then continued data collection (in sets of approximately 8-12 participants) until Bayesian statistics converged on support for or against the null hypothesis. Unlike frequentist statistics when the result is null, Bayesian statistics are consistent, and will converge on the true model with more data (Rouder et al., 2009).

3.2.4 Results and discussion

As in Experiment 1, our first hypothesis was that participants might be invariant to differences between instruments. Contrary to the predictions of strict invariance, and analogous to the result of Experiment 1, we found that discrimination was worse when participants compared notes from different instruments than when they compared notes from the same instruments (main effect of same vs. different instruments, $F(1,85)=140.34$, $p<.0001$, $\eta_p^2=.62$, Figure 2c). Participants were nonetheless well above chance at discriminating notes from different instruments (performance was above chance even for the smallest pitch change tested, .11 semitones, t-test against .5, $t(85)=10.04$, $p<.0001$), again demonstrating some degree of spectral invariance.

Our second hypothesis was that harmonicity would be critical to spectral invariance, such that there would be a harmonic advantage in the different instrument condition but not in the same instrument condition (i.e., an interaction). There was a statistically significant interaction between Harmonicity and Same vs. Different instruments ($F(1,85)=11.60$, $p=.001$, $\eta_p^2=.12$). However, the effect was small, and in the opposite direction from the prediction of the hypothesis, being driven by better performance for Inharmonic notes in the Same Instrument condition. As with vowels in Experiment 1, there was no significant effect of Harmonicity with Different Instruments for the smallest three step sizes ($F(1,85)=0.002$, $p=.96$, $\eta_p^2<.0001$), and evidence for the null hypothesis (Bayes factor BF_{incl} , specifying a multivariate Cauchy prior on the effects, was .07, providing strong support of the null hypothesis (JASP, 2020)). There was a significant effect of Harmonicity for Same Instruments with small step sizes ($F(1,85)=15.66$, $p=.0002$, $\eta_p^2=.16$), but it was due to performance being marginally better for Inharmonic notes, for reasons that are unclear. The results are thus again overall inconsistent with the hypothesis that harmonicity underlies spectral invariance.

As with vowels in Experiment 1, discrimination was worse for Inharmonic notes when the step size was large. At a 9 semitone difference there was again a highly significant effect of Harmonicity ($F(1,85)=136.71$, $p<.0001$, $\eta_p^2=.62$; this drove an overall effect of Harmonicity across all step sizes, $F(1,85)=22.16$, $p<.0001$, $\eta_p^2=.21$). This again produced a highly significant interaction between step size and Harmonicity ($F(1,85)=48.93$, $p<.0001$, $\eta_p^2=.37$).

All together, these results suggest that pitch discrimination is not completely invariant to naturally occurring differences in spectral shape – performance was worse overall when participants compared notes played on different vs. the same instrument. Second, at small pitch differences, performance was similar for harmonic and inharmonic conditions when listeners compared notes played on different instruments, suggesting that any invariance participants exhibit in discriminating different instruments is not supported by representations of f_0 . However, for larger pitch changes, we observed a deficit for inharmonic notes for both Same and Different Instrument conditions. This could indicate that when listening to natural sounds like instruments, listeners rely on spectral cues to register small pitch changes (yielding the matched performance for Harmonic and Inharmonic conditions), but transition to using the f_0 for large changes (better performance for Harmonic conditions).

3.2.5 Post Hoc Spectral Centroid Analysis

It seemed plausible that the deficit we observed in the Different Instrument condition could depend on whether the difference in the spectral center-of-mass (centroid) between the two instrument notes was congruent or incongruent with the pitch change. To investigate this issue, we performed an additional analysis in which we calculated the spectral centroid of each note and then separated Different Instrument trials into two sub-sets of trials: ‘Congruent’ and ‘Incongruent’. Trials were categorized as Congruent if the f_0 and spectral centroid shifted in the same direction. For instance, if the f_0 was higher in the second note than the first, so was the spectral centroid. Trials were categorized as Incongruent if the f_0 and spectral centroid moved in opposite directions.

To determine the spectral centroid of instruments we first ran instrument notes through STRAIGHT and resynthesized them to be pitch-flattened and harmonic (to match the stimuli used in the

experiment), and normalized them to have an rms level of .05. We then estimated the power spectral density using Welch's method. We used the default parameters of the MATLAB implementation of this method to obtain an estimate of the power at each frequency (1422 samples per window, 50% overlap, defined using Hamming windows with 42.5 dB sidelobe attenuation). We expressed the frequencies in Equivalent Rectangular Bandwidths (ERBs) using the formula of Glasberg and Moore (Glasberg & Moore, 1990), then to estimate the spectral centroid we took a weighted average of these frequencies. The weights were the power (in dB) of each frequency relative to the noise floor, which we estimated by eye to be -65 dB, with values below the noise floor set to 0 so that frequencies with low power would not contribute to the weighting. We repeated this procedure for 12 different notes (covering a full octave) for each instrument and averaged together the 12 estimates to obtain an average spectral centroid for each instrument: (from lowest to highest) Ukulele, Pipe organ, Cello, Oboe, and Baritone Saxophone.

For each of the 20 pre-generated stimulus sets described above, we classified each trial as either Congruent or Incongruent, and then separately analyzed the two sets of trials. Since the Different instrument pairs were randomly selected on each trial (because the experiment was not originally designed to accommodate the congruency analysis), the Congruent and Incongruent trials were not exactly balanced (overall, 51.4% of trials were classified as Congruent, and 48.6% of trials were classified as Incongruent).

As observed in the inset to Figure 2c, performance in the Incongruent condition was significantly worse than that in the Congruent condition ($F(1,85)=199.63$, $p<.0001$, $\eta_p^2=.70$). This analysis suggests that the spectral shape of natural sounds biases pitch discrimination judgments.

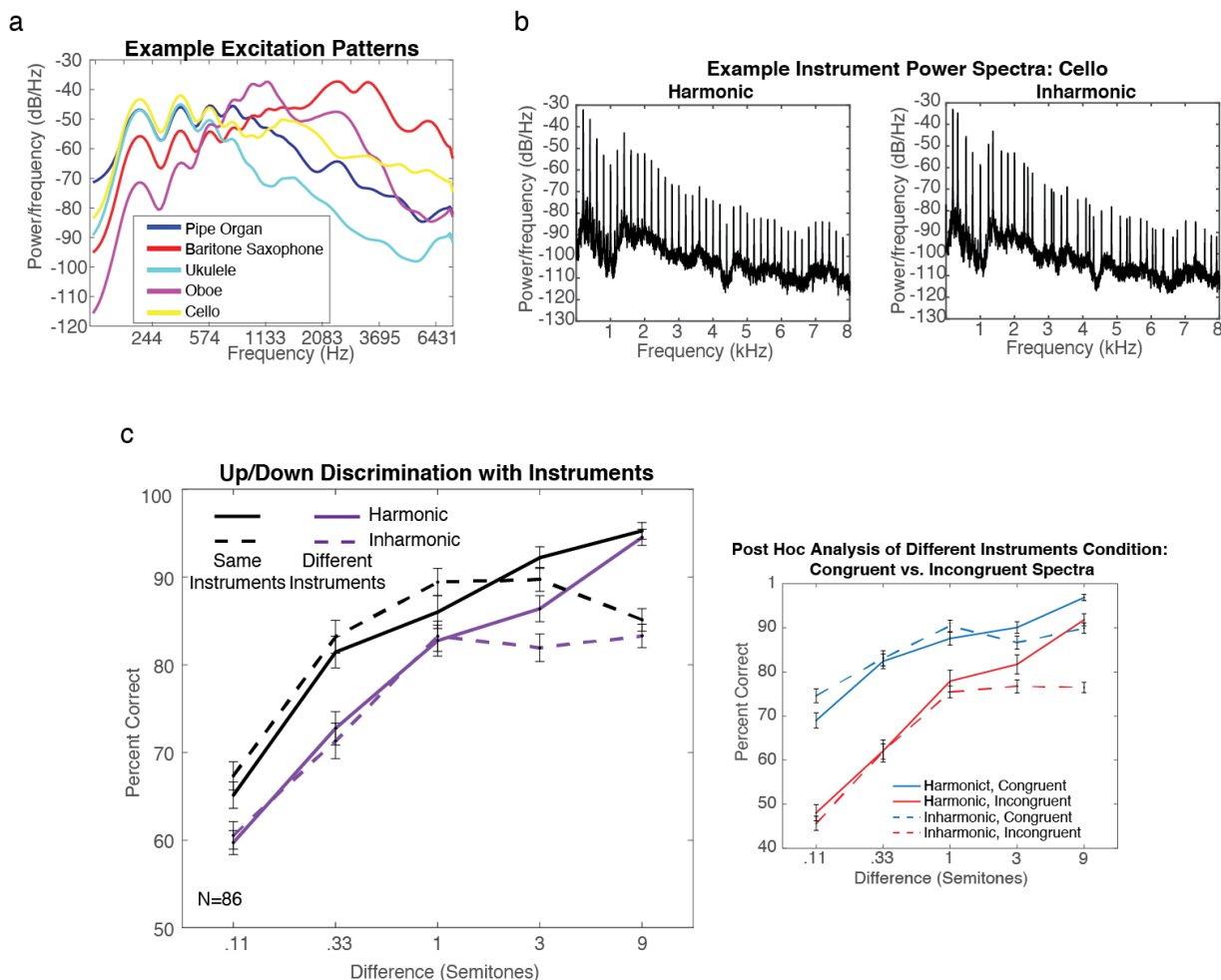


Figure 2: Stimuli and results for Experiment 2, Pitch discrimination with notes played on the same vs. different instruments

a. Example excitation patterns for notes with the same f_0 (200 Hz) played on the five instruments used in Experiment 2. Excitation patterns contain peaks for individual harmonics, but also reveal the differences in spectral shape between the instruments. b. Example power spectra of harmonic and inharmonic instrument notes used in Experiment 2. Note the irregularly spaced frequencies in the inharmonic notes, due to the frequency jittering. c. Results from Experiment 2. Error bars denote standard error of the mean. Left panel shows Same and Different Instrument (main results). Right panel shows the results of a post-hoc analysis in which ‘Different Instrument’ trials were grouped by whether the spectral centroid of the notes moved congruently or incongruently with the f_0 shift (‘Congruent’ and ‘Incongruent’ subsets of trials).

3.3 Effects of Musicianship

It was a priori unclear whether any of the effects we were testing in Experiments 1 and 2 could be influenced by musical training. Western musical training has been associated with lower pitch discrimination thresholds (Bianchi et al., 2016; Kishon-Rabin et al., 2001; McDermott et al., 2010; McPherson & McDermott, 2018; Micheyl et al., 2006; Spiegel & Watson, 1984). It seemed plausible that musicians might have more experience comparing notes across different instruments than non-musicians, and thus might be more robust to spectral differences. However, we found no evidence for

such differences in robustness. While there was a main effect of musicianship in both experiments (two-way ANOVAs, Experiment 1: $F(1,80)=16.16$, $p=.0001$, $\eta_p^2=.17$ Experiment 2: $F(1,84)=8.43$ $p=.0047$, $\eta_p^2=.09$), driven by better discrimination performance in musicians, we did not observe significant interactions between musicianship and same/different vowels or instruments (mixed model ANOVAs, Experiment 1, Vowels: $F(1,80)=1.93$ $p=.18$, $\eta_p^2=.02$, Experiment 2, Instruments: $F(1,84)=0.51$ $p=.47$, $\eta_p^2=.006$). There was likewise no interaction between musicianship and harmonicity (mixed model ANOVAs, Experiment 1, Vowels: $F(1,80)=1.88$, $p=.18$, $\eta_p^2=.02$), Experiment 2, Instruments: $F(1,84)=0.88$, $p=.35$, $\eta_p^2=.01$). These findings suggest that the effects of spectral variation on pitch perception are not strongly dependent on formal musical training. Given the lack of a musicianship effect in these two experiments, we did not analyze musicianship in subsequent experiments.

3.4 Experiment 3: Discrimination of synthetic tones with extreme spectral differences

3.4.1 Purpose and procedure

Experiment 3 was designed to further explore the biasing effects observed in Experiment 2. We aimed to determine whether the bias would remain in a setting where listeners were forced to rely on representations of f_0 to complete the task (Figure 3a). To address this question, we tested participants with synthetic tones whose harmonics were completely non-overlapping from note-to-note (as if the result of unnaturally steep rectangular filters that completely attenuate harmonics outside their passband). The filters that produce most natural sounds, as in the instrument tones in Experiment 2, tend to modestly attenuate harmonics rather than render them inaudible. This modest attenuation leaves harmonics in common between notes, and our results suggest that listeners can compare these frequencies to detect pitch differences even when notes are inharmonic. However, when the two tones contain entirely separate subsets of the harmonic series, discrimination should be impossible when tones are altered to be inharmonic, as there are no common frequency components to compare. Above-chance performance when tones are harmonic thus requires a representation of the f_0 . We asked whether listeners are biased by the spectrum in such conditions, by testing listeners using notes whose spectral content differed to be either 'Congruent' or 'Incongruent' with the f_0 shift, comparable to Experiment 2.

As in Experiments 1 and 2, participants heard two tones and were asked whether the second tone was higher or lower than the first tone. Participants completed 30 trials per condition for all Same Harmonic conditions. We had originally intended to average together 'Congruent' and 'Incongruent' trials in a 'Different Harmonics' condition, thus participants only completed 15 Congruent trials and 15 Incongruent trials. Congruent and Incongruent conditions were randomly intermixed (randomized independently for each participant).

PsychoToolbox for MATLAB (Kleiner et al., 2007) was used to play sound waveforms. Sounds were presented to participants at 70 dB SPL over Sennheiser HD280 headphones (circumaural) in a soundproof booth (Industrial Acoustics). Participants were consented before beginning the experiment and were paid a fixed hourly rate for their participation. All participants self-reported normal hearing.

3.4.2 Stimuli

The two tones on each trial contained sets of equal-amplitude harmonics, added in sine phase. For 'Congruent' and 'Incongruent' conditions, one tone contained harmonics 1-5, or 1-6 (low harmonics), and the other note contained harmonics 6-12 or 7-12, respectively (high harmonics) (Figure 3a). In Congruent conditions if the second tone's f_0 (or nominal f_0 , for inharmonic conditions) was higher than the first, the low harmonics were used in the first note and the high harmonics were used in the second note, and vice versa when the second note's f_0 was lower than the first. Incongruent conditions had the inverse relationship between the f_0 and the spectra. For 'Same Harmonics' conditions, the same set of high or low harmonics (equiprobably) were used in each tone. The tones were 400 ms in duration, windowed with 20 ms half-Hanning windows, and separated by a 200 ms silent interval.

For all conditions, notes could be either harmonic or inharmonic. All inharmonic tones were made

inharmonic as described above in *Inharmonic Stimuli*. The f_0 (or nominal f_0 , for inharmonic conditions) of the first tone of each trial was randomly selected from a log uniform distribution spanning 200 to 400 Hz. The f_0 of the second tone was either higher or lower (equiprobably) than that of the first tone by the necessary step size.

3.4.3 Participants

Experiment 3 was completed in lab. Trials for Experiment 3 were intermixed with those from another experiment in which tones were embedded in background noise. The results of that other experiment are not reported here. All participants were native English speakers residing in the Greater Boston area. 27 participants completed the experiment but 3 performed below 55% correct across all conditions, and were excluded from further analysis. The data presented here is from the remaining 24 participants (11 self-identified as female, 13 as male, 0 as non-binary; mean age = 32 years, S.D.=14 years). 13 participants had over 4 years of musical training (self-reported, mean = 17, S.D.=13). This sample size allows a 90% chance of seeing a small-to-moderate effect size (Cohen's $f = .2$) at a $p < .05$ level of significance.

3.4.4 Results and discussion

Although performance was above chance in all conditions, we found that the impairment observed in Experiment 2 persisted for synthetic tones with non-overlapping harmonics (Figure 3b, left). Performance was better when listeners heard the same harmonics from note-to-note than in both Congruent and Incongruent conditions (Same Harmonics vs. Congruent, Harmonic tones, $F(1,23)=7.26$, $p=0.013$, $\eta_p^2=.24$; Same Harmonics vs. Incongruent, Harmonic tones, $F(1,23)=21.51$, $p=0.0001$, $\eta_p^2=.48$). While there was not the significant main effect of congruency that would signify a spectral bias ($F(1,23)=2.06$, $p=0.16$, $\eta_p^2=.08$), this appears to be due to near-chance performance for the smaller step sizes. At the 1 semitone step size, performance in the Incongruent condition was significantly worse than that for the Congruent condition ($t(23)=2.15$, $p=.005$). The spectral bias observed in Experiment 2 thus appears present in these conditions as well.

As in previous experiments with synthetic tones (Faulkner, 1985; McPherson & McDermott, 2018; McPherson & McDermott, 2020; Micheyl et al., 2010; Moore & Glasberg, 1990), and consistent with Experiment 1, performance was similar for Harmonic and Inharmonic tones when the two tones had the same harmonics (Figure 3b, black solid and dashed lines, $F(1,23)=0.65$, $p=0.43$, $\eta_p^2=.03$). However, as expected, listeners could not perform the task with Inharmonic tones whose harmonics did not overlap, with performance entirely driven by the direction of spectral change: there was no main effect of step size ($F(1,23)=1.00$, $p=.38$, $\eta_p^2=.04$). This suggests that – as intended based on the design of the tones – listeners rely on the f_0 when notes are harmonic and spectral differences are sufficiently extreme. But even in this case, spectral invariance is incomplete.

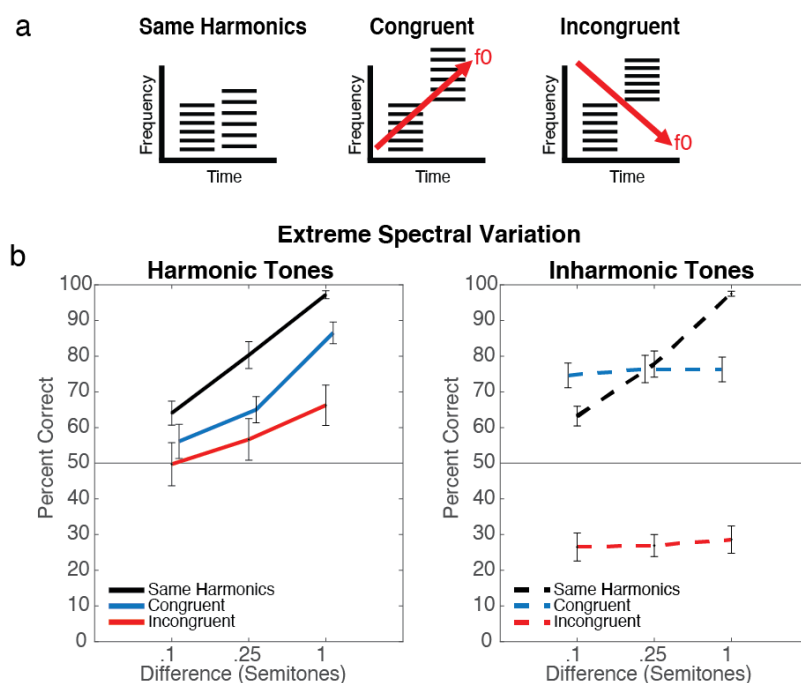


Figure 3: Stimuli and results for Experiment 3, Discrimination of synthetic tones with extreme spectral differences

a. Stimulus configurations. The two tones on a trial could either have the same set of harmonics, or non-overlapping sets of harmonics, arranged so that the change in spectral centroid was either Congruent or Incongruent with the change in f_0 . b. Results from Experiment 3, grouped by harmonicity. Error bars denote standard error of the mean.

3.5 Experiment 4: Effect of time delay on spectral invariance

3.5.1 Purpose and procedure

Experiments 2-3 showed that pitch discrimination judgments are biased by the spectral content of notes, even in conditions where listeners must rely on the f_0 to perform the task because the harmonics are completely non-overlapping. There are several alternative hypotheses that could explain this bias. One possible explanation is that the representation of a sound's f_0 is biased by the sound's spectrum. Another possibility is that the f_0 representation is itself unbiased, but that listeners are unable to base decisions exclusively on this representation, with spectral differences between notes biasing their judgments.

In Experiment 4 we attempted to disambiguate these hypotheses using time delays. In previous work we found that listeners become more reliant on representations of the f_0 when making judgments about pitch across a delay (McPherson & McDermott, 2020), as though the f_0 representation is better remembered than representations of the spectrum. If the f_0 representation is itself unbiased, with spectral biases coming from suboptimal use of cues at a decision stage, it seemed possible that the bias would be reduced over a delay.

As in Experiment 3, tones either contained the same set of harmonics or non-overlapping harmonics. On the non-overlapping trials, the spectrum could be either Congruent or Incongruent with the f_0 change (Figure 4a). Listeners heard two notes during each trial, either played back-to-back or separated by a 3 second delay (Figure 4b). We only tested participants using Harmonic notes. The task was always to say whether the second tone was higher or lower than the first, and participants received feedback (Correct/Incorrect) after each trial. Participants completed 24 trials per condition.

3.5.2 Stimuli

Tones were identical to those used in Experiment 3, except they were either played back-to-back, or separated by a 3 second delay. The f_0 of the first tone of each trial was randomly selected from a log uniform distribution spanning 200 to 400Hz.

3.5.3 Participants

194 participants passed the headphone check and completed Experiment 4 online. 87 had overall performance (averaged across all conditions) below 55% correct and were removed from further analysis. The large number of participants excluded based on this criterion partly reflects the fact that data was collected during a time when Amazon's Mechanical Turk platform was experiencing a higher-than-normal level of fraud, and we observed that many participants were performing at chance levels. Of the 107 remaining participants, 39 identified as female, 67 as male, 1 as non-binary. The average age of these participants was 37.2 years (S.D.=10.0). 35 participants had four or more years of musical training, with an average of 10.6 years (S.D.=10.1).

To determine sample size, we performed a power analysis by bootstrapping pilot data (from an earlier version of Experiment 4 in which listeners heard all but the 9-semitone step size condition). For each possible sample size, we computed bootstrap distributions of the interaction F statistic (Congruent/Incongruent vs. No Delay/Delay). We found that a sample size of 11 yielded a 95% chance of seeing the interaction present in our pilot data at $p < .01$ significance level. However, based on our pilot data we also wanted to have sufficient power to detect an effect of delay specifically on the Incongruent condition (where we found that performance improved with delay). A power analysis using G*Power (Faul et al., 2007) indicated that we needed at least 68 participants to be 95% sure of seeing this an effect half of that observed in the pilot study ($\eta_p^2 = .13$) at a $p = .01$ significance level. In practice, we ran the experiment in small batches and ended up with more than this number of participants.

3.5.4 Results and discussion

As shown in Figure 4b, the effects of the spectrum decreased with delay, producing a significant interaction between Congruent/Incongruent conditions and Delay ($F(1,106) = 21.89$, $p < .0001$, $\eta_p^2 = .17$). We replicated the pronounced effect of spectral variation observed in Experiment 3 when there was no delay (performance was worse for Congruent than for Same Harmonics ($F(1,106) = 8.13$, $p = .005$, $\eta_p^2 = .07$), and worse for Incongruent than Congruent ($F(1,106) = 48.95$, $p < .0001$, $\eta_p^2 = .32$). However, these effects were reduced when there was a delay between tones, due to the delay having different effects in different stimulus conditions. The delay reduced performance for the Same Harmonics condition (Figure 4c, significant main effect of delay; $F(1,106) = 17.48$, $p < .0001$, $\eta_p^2 = .14$). But there was no significant effect of delay in the Congruent Spectra condition ($F(1,106) = 2.48$, $p = .12$, $\eta_p^2 = .02$). And the delay in fact *improved* performance in the Incongruent condition (significant effect of delay, but in the opposite direction; $F(1,106) = 15.62$, $p = .0001$, $\eta_p^2 = .13$).

It is not obvious how to account for these results without positing that the representations of pitch and spectral shape are separate. If the representation of pitch were itself biased, the bias should persist across delays. In particular, the delay should impair performance for both Congruent and Incongruent conditions. In this respect the most diagnostic result is the counterintuitive improvement with delay in the Incongruent condition. The change in bias with delay suggests that the representations of pitch and spectral shape are independent, and decay at different rates, with the effect of the spectrum on pitch judgments decreasing over time.

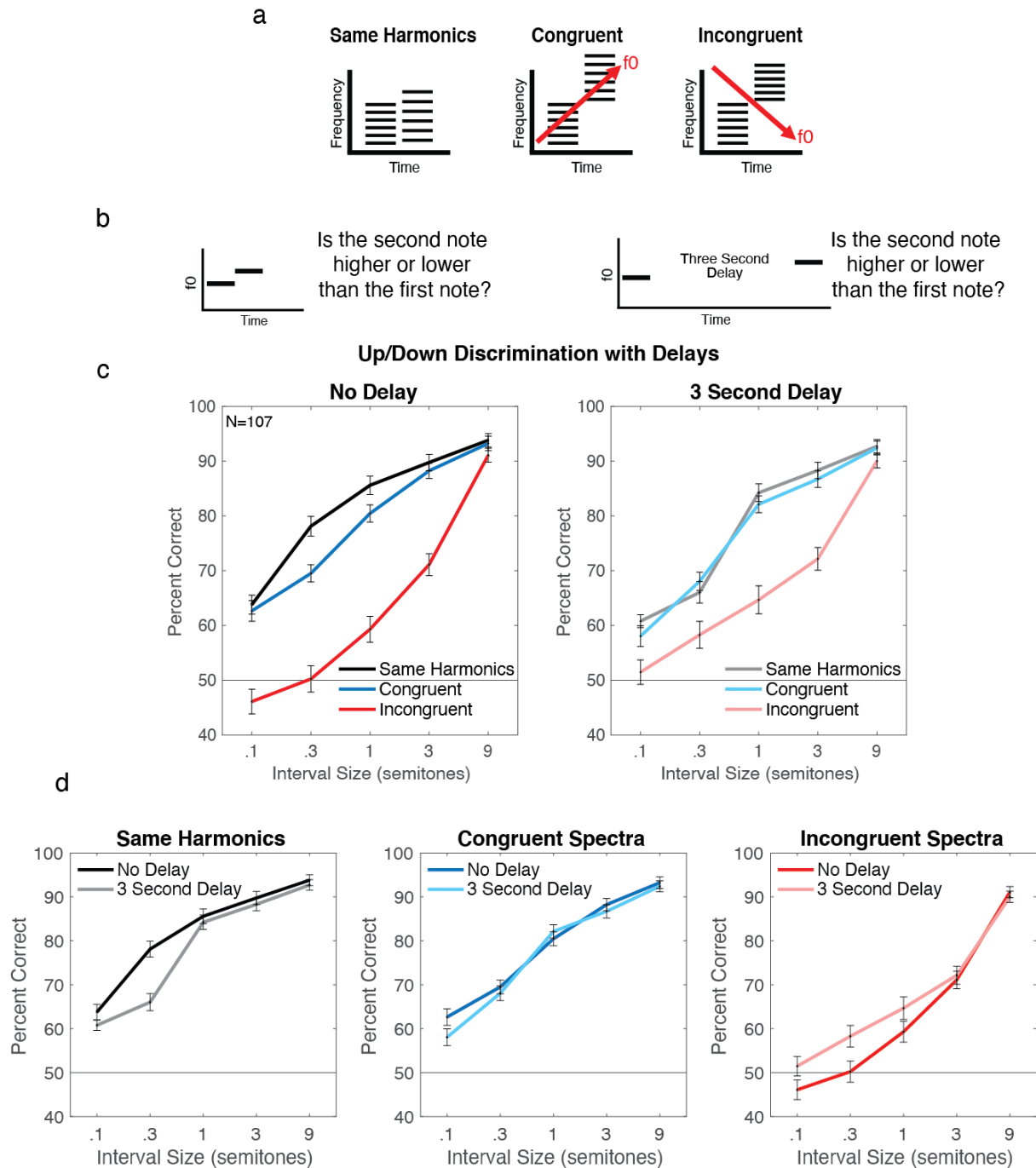


Figure 4: Stimuli and Results for Experiment 4, Effect of delays on discrimination performance with synthetic tones

a. Stimulus configurations. The two tones on a trial could either have the same set of harmonics, or non-overlapping sets of harmonics, arranged so that the change in spectral centroid was either Congruent or Incongruent with the change in f_0 . Tones could either be presented back-to-back (no delay; top row), or separated by a three second delay (bottom row) b. Results from Experiment 4, grouped by delay conditions. Here and in c, error bars denote standard error of the mean. c. Results from Experiment 4, grouped by spectra condition.

3.6 Experiment 5: Spectral invariance of musical interval perception

3.6.1 Purpose and procedure

Experiment 5 was intended to test whether the spectral bias we observed across several basic pitch discrimination tasks (Experiments 2-4) would generalize to judgments about the magnitude of pitch differences (intervals). This question was previously addressed by Russo and Thompson, who asked listeners to rate interval sizes on a scale of 1 to 5, and found that timbre changes could augment or diminish the rated interval size (Russo & Thompson, 2005). The purpose of our experiment was to test whether this effect would occur for a musical judgment based on pitch intervals. In our experiment, listeners heard two notes and were asked whether the notes matched the starting interval in ‘Happy Birthday’, specifically, the two-semitone difference between the last syllable of ‘Happy’ and the first syllable of ‘Birthday’ (Figure 5a). Listeners either heard notes separated by the correct interval of 2 semitones, or by an interval that was mistuned by up to a semitone in either direction (2 semitones \pm .5 or 1 semitones). Participants completed 6 trials per condition.

We tested four conditions: Harmonic, Inharmonic, Consistent, and Inconsistent (Figure 5b). The Harmonic and Inharmonic conditions presented pairs of tones identical to those in Experiment 3 apart from the different pitch intervals used. The purpose of these conditions was to establish baseline performance on the task, and to help characterize the representations used for the task (namely, whether performance was dependent on representations of the f_0). The latter seemed advisable given that this was not a task we had previously used.

The ‘Consistent’ and ‘Inconsistent’ conditions were analogous to the ‘Congruent’ and ‘Incongruent’ conditions in Experiments 3 and 4. However, the second note in each trial was always higher than the first note, so instead of the change in spectral centroid being aligned or mis-aligned with the f_0 change, it was aligned or mis-aligned with the positive or negative mistuning from the correct two-semitone interval. For example, in the ‘Consistent’ condition, if the pitch mistuning was positive (the presented interval was larger than two semitones), then the first note contained the lower set of harmonics and the second note contained the higher set of harmonics, and vice versa for a negative mistuning. Our prediction was that the spectral shift might increase or decrease a listener’s estimate of the interval size, making the correct answer more apparent in the ‘Consistent’ case, and less apparent in the ‘Inconsistent’ case.

3.6.2 Stimuli

On each trial participants were presented with two synthetic complex tones in succession (with no delay between tones). Tones were identical to those used in Experiment 3. The f_0 (or nominal f_0 , for Inharmonic conditions) of the first tone of each trial was randomly selected from a log-uniform distribution spanning 200 to 400 Hz. The f_0 of the second tone was either 2 semitones above the first (to match the first interval in ‘Happy Birthday’ or mistuned from 2 semitones by \pm .5 or 1 semitone. On Consistent and Inconsistent trials the spectral centroid of the notes moved in the same or opposite direction as the mistuning, as described above. There were an equal number of ‘Happy Birthday’ and ‘Not Happy Birthday’ trials throughout the experiment.

3.6.3 Participants

252 participants passed the headphone check and completed Experiment 3 online. 203 had performance (averaged across all conditions) below a d' of .67 and were removed from further analysis. We chose this exclusion criteria based on a small pilot study completed in the lab before the COVID-19 pandemic, with 5 participants completing the task with Harmonic and Inharmonic tones. Across harmonicity conditions, the average d' across the two smallest mistunings (\pm .5 semitones) was .67, so we used this as a conservative estimate of good performance among online participants. As in Experiment 4, the large number of participants excluded based on this criterion partly reflects the fact that data was collected during a period when Amazon’s Mechanical Turk platform was experiencing a

higher-than-normal level of fraud, and we observed that many participants were performing at chance levels. Of the remaining 49 participants, 23 identified as female, 24 as male, 2 as non-binary. The average age of these participants was 38.5 years (S.D.=11.5). 23 participants had four or more years of musical training, with an average of 14.6 years (S.D.=11.7).

The target sample size was determined based on a separate pilot experiment, run online, where participants only heard Harmonic and Inharmonic conditions. In this experiment participants heard 10 trials per condition, rather than 6. To estimate an effect size with fewer trials per participant we bootstrapped from pilot data, sub-sampling 6 trials from each participant per condition and calculated the Harmonic vs. Inharmonic effect size 10,000 times. The average effect size was $\eta_p^2=.23$. We predicted that effects of the spectrum might be smaller than effects of harmonicity. A power analysis indicated that we would need 38 participants to observe an effect 1/2 the size of the effect of harmonicity at a p value of .01, 95% of the time (Faul et al., 2007).

3.6.4 Results and discussion

Participants were well above chance in the Harmonic condition, with performance improving with the magnitude of the mistuning, as expected. Performance remained above chance in the Inharmonic condition, but was substantially worse than in the Harmonic condition (Figure 5c, $F(1,48)=17.86$, $p=.0001$, $\eta_p^2=.27$). These results complement previous findings that tasks related to musical interval perception are impaired with inharmonic tones, thus likely relying on f0-based pitch (McPherson & McDermott, 2018). The above-chance performance with inharmonic tones suggests that listeners can nonetheless extract some information about interval size from shifts in the spectrum, potentially by attending to the lowest frequency component of the tones.

Parallel to the spectral biases found for basic discrimination in Experiments 2-4, performance was better in the Consistent condition compared to the Inconsistent condition ($F(1,48)$, 6.94, $p=.011$, $\eta_p^2=.13$). This suggests that the limited invariance of pitch judgments is not specific to up-down pitch discrimination. This result is also consistent with the idea that the bias in pitch judgments originates in representations of relative pitch. Unlike a basic discrimination task, here the decision involves judging whether the heard interval is the same or different as a remembered interval. It is not obvious how to explain the effect of spectral shape on the results without positing that the note-to-note change in the spectrum influences a representation of the pitch interval (note-to-note change in pitch). Relative pitch representations may thus be biased by changes in the spectrum, even if the f0 representation they are based on is not. These biases in relative pitch could explain the effects of spectral shape seen on basic pitch discrimination (Experiments 1-4).

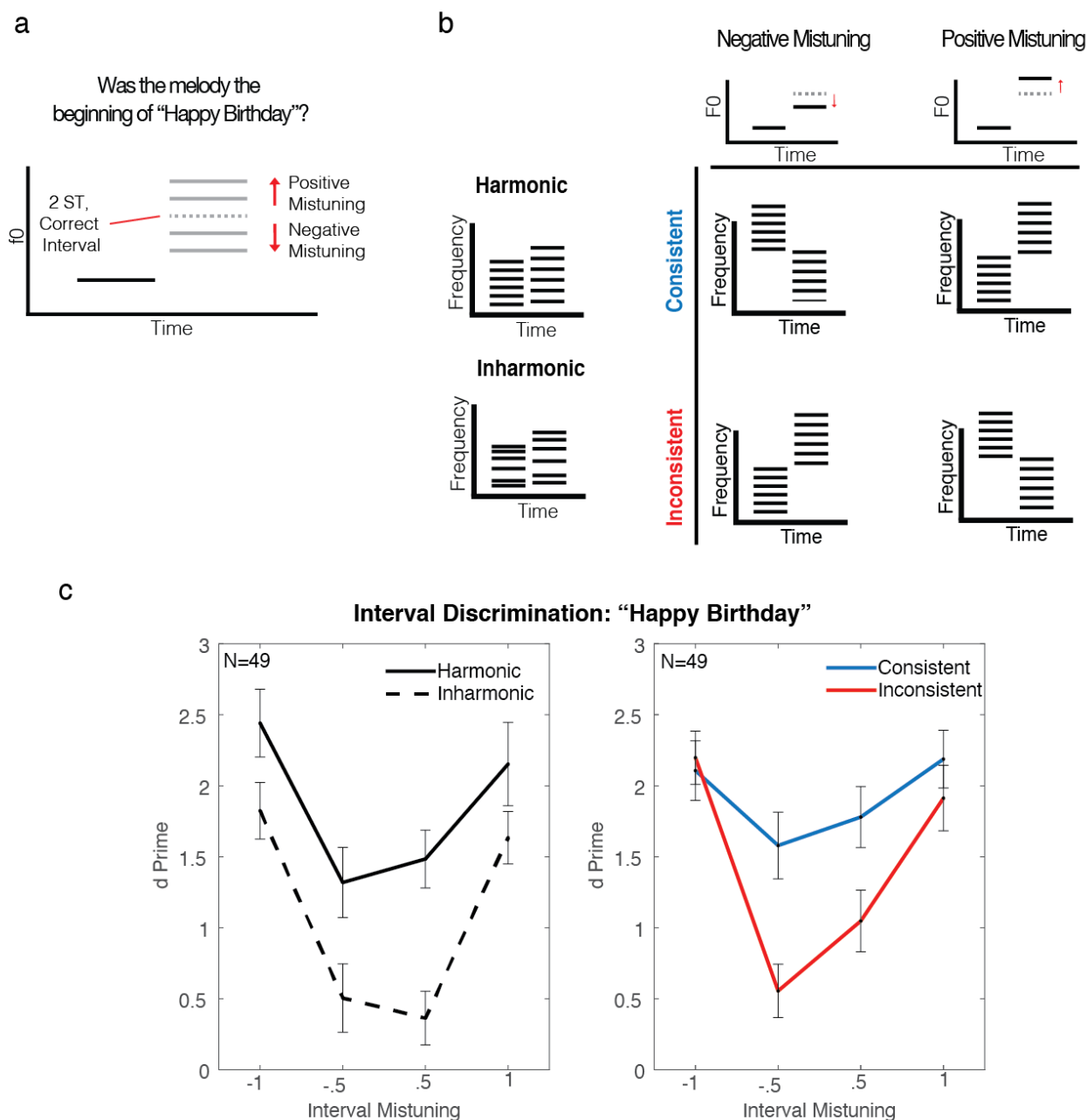


Figure 5: Stimuli and Results for Experiment 5, Discrimination of intervals composed of synthetic tones with extreme spectral differences

a. Schematics of a trial. b. Schematics of stimulus conditions. Two right-most columns show the relationship between the spectra and f_0 in ‘Consistent’ and ‘Inconsistent’ conditions. c. Results from Experiment 4. Error bars denote standard error of the mean.

4 GENERAL DISCUSSION

We tested the extent to which listeners can discriminate sounds differing in spectral shape, using both natural sounds such as instruments and voices, and synthetic tones with extreme levels of spectral variation. Our first question was whether pitch discrimination is invariant to real-world spectral differences across sounds. We found that listeners were able to discriminate f_0 differences between notes of different instruments, or between different vowels, but their performance was worse than when comparing the same instrument or vowel, being biased by changes in spectral centroid. These results suggest that listeners are only partially invariant to differences in spectral shape between natural sounds. This indicates that previous findings of imperfect invariance with synthetic tones generalize beyond stimuli with unnatural spectral differences. Our second question was whether any spectral invariance in pitch perception is mediated by representations of f_0 , in which case discrimination should be worse for inharmonic sounds that differed in spectral shape. We observed inharmonic deficits when f_0 differences between otherwise natural sounds were sufficiently large, but these deficits were similar for conditions where the spectral shape (vowel or instrument) was the same vs. different. The only cases in which spectral variation produced impairments for inharmonic stimuli involved synthetic tones with extreme spectral differences between notes. Our third question was whether the limits of invariance reflect biases in representations of the f_0 vs. processes that operate on those representations to mediate decisions. We found that the spectral bias of pitch judgments decreased across a delay, suggesting that the bias is not exerted on the representation of a sound's absolute pitch. All together, our results show that human pitch perception exhibits some degree of limited spectral invariance, enabling pitch comparisons between natural sounds, but that this invariance does not normally require representations of f_0 , or make use of the fact that f_0 representations are apparently unbiased.

4.1 Real-world discrimination is somewhat robust to timbre, independent of f_0 -based pitch

We built on previous studies examining pitch discrimination between tones differing in spectral shape (Allen & Oxenham, 2014; Micheyl & Oxenham, 2004; Moore et al., 1992; Moore & Glasberg, 1990; Russo & Thompson, 2005; Singh & Hirsh, 1992; Warrier & Zatorre, 2002) by using natural sounds that varied in their spectra due to differences in filtering from the vocal tract or instrument body. In both Experiment 1 (speech) and Experiment 2 (musical instruments), participants showed above-chance up-down discrimination irrespective of whether the instruments or speakers being compared were the same or different. This result confirms that pitch perception is somewhat robust to everyday spectral differences between sounds. However, performance was nonetheless substantially worse when different instruments or vowels were compared. Critically, the impairments resulting from between-sound spectral variation were no greater for inharmonic sounds. For small f_0 differences (3 semitones and lower) discrimination performance was similar for harmonic and inharmonic conditions. For larger f_0 differences discrimination was impaired for inharmonic tones, but these effects held irrespective of whether the instruments or vowels were the same or different. This result suggests that human robustness to everyday spectral variation does not rely on representations of the f_0 . The results are consistent with other evidence that listeners do not normally use f_0 representations when discriminating sounds presented back-to-back with modest f_0 differences (McPherson & McDermott, 2018; McPherson & McDermott, 2020), but suggest that this is the case even for spectral differences between sounds typically encountered in natural conditions.

4.2 Invariance to extreme spectral differences is mediated by f_0 -based pitch

In Experiment 3, we introduced extreme (and unnatural) spectral changes between notes by presenting nonoverlapping subsets of frequency components. In these conditions, listeners were only able to make accurate discrimination judgments when the stimuli were harmonic, indicating a reliance on f_0 -based pitch. However, performance with harmonic tones was nonetheless worse when the two notes did not share harmonics. Both findings are consistent with previous studies examining pitch discrimination

between sounds that vary in their spectra (Allen & Oxenham, 2014; Micheyl & Oxenham, 2004; Moore et al., 1992; Moore & Glasberg, 1990; Singh & Hirsh, 1992; Warrier & Zatorre, 2002). It is thus clear that listeners can use representations of the f_0 when there is extreme spectral variation between sounds, but this paper provides evidence that such experimental manipulations are not fully representative of what happens when we discriminate natural sounds. Evidently the spectral differences between natural sounds are sufficiently modest that any invariance does not require f_0 -based pitch.

4.3 Necessity of f_0 -based pitch for registering large pitch changes

In Experiments 1 and 2 we found that when pitch changes were sufficiently large, performance with inharmonic stimuli was impaired relative to harmonic stimuli. This impairment plausibly reflects ambiguities in the correspondence between frequency components due to the filter that is present for most natural sounds. When pitch changes are small the correspondence between frequency components of two sounds being compared is likely to be unambiguous. This is because the f_0 difference is small relative to the harmonic spacing, such that a resolved frequency component in the first sound has a close match at the corresponding harmonic in the second sound. However, when pitch changes become sufficiently large, the correspondence between components may become ambiguous when a filter attenuates the upper and lower ends of the source spectrum (Figure 1f). f_0 -based pitch evidently helps listeners to register such large changes in f_0 . In a previous study we observed analogous effects with synthetic tones passed through a fixed bandpass filter (replotted in Supplementary Figure 1 for comparison). The contribution of the present experiments is to show that such effects occur in natural sounds and thus have relevance for real-world listening. The inharmonic deficit at large intervals was present for both speech and music sounds, providing some evidence for a domain-general phenomenon.

4.4 Origins of spectral biases in pitch judgments

The results of Experiments 2-3 left it unclear whether pitch discrimination biases from spectral variation are the result of bias in the pitch estimation process or bias at some subsequent stage, such as the task decision. We leveraged the effects of time delay on pitch discrimination to explore these possibilities. Previous experiments have suggested that listeners become more reliant on representations of f_0 when comparing sounds across a time delay (McPherson & McDermott, 2020), potentially because representations of the spectrum deteriorate more rapidly over time than those of the f_0 . This result raised the possibility that time delays might help to distinguish the origin of the spectral bias. We reasoned that if the bias occurs at a stage where the pitch of different sounds is being compared, and if spectral shape information decays more over time than does information about the f_0 , then the observed bias might change with the delay. Experiment 4 substantiated this prediction – the bias in discrimination judgments decreased across a delay, and in particular, performance improved with an inter-note delay for trials where the spectral shape and f_0 were incongruent. These results thus suggest that the bias in pitch judgments is likely not the result of bias in the representation of the pitch of individual sounds, but rather arises at a subsequent comparison stage. This conclusion is consistent with findings that pitch matching is not much affected by note timbre (Russo & Thompson, 2005).

Experiment 5 raised the possibility that spectral biases occur at the stage of relative pitch representations. The task in that experiment required listeners to compare the stimulus to their memory of the first interval in ‘Happy Birthday’, and make a same-different judgment. The results showed evidence of timbre differences biases on pitch intervals, replicating findings of Russo and Thompson (2005) but with a musical judgment. In contrast to the biases seen in basic pitch discrimination (Experiments 2-4), it is not obvious how to explain the bias in interval judgments purely with a decision stage. With up-down pitch discrimination judgments, the bias in decisions could result from a competing direction signal from the spectral shape being difficult for listeners to ignore when they choose “up” or “down”. In principle, such decision-stage interference could also occur with the interval size ratings measured by Russo and Thompson. But in the interval task we used in Experiment 5, the judgment was essentially a

same/different judgment comparing the stimulus to a memory representation. Moreover, the pitch interval stimulus always had the same direction, and differed only in its magnitude. It is thus less obvious how the direction of the timbral change from note to note could aid or impair the decision process. One possibility is that changes in spectral shape are integrated with changes in pitch at the stage of relative pitch representations (Russo & Thompson, 2005), which in our experiment could result in augmented or diminished representations of the pitch interval between the two notes. Biases in relative pitch could also account for the biases seen in basic pitch discrimination, on the assumption that the decision variable is a representation of relative pitch. The influence of coarse spectral changes on relative pitch is consistent with findings that humans extract melody-like representations from such changes (Cousineau et al., 2014; Graves et al., 2014; Graves et al., 2019; McDermott et al., 2008; Siedenburg, 2018).

4.5 Limitations

Our natural sound stimuli were constrained by available corpora. We attempted to test the upper end of naturally occurring spectral variation by choosing pairs of vowels (Experiment 1) and instruments (Experiment 2) with maximally different excitation patterns. We think it likely that the extent of variation in our stimuli is representative of that encountered in everyday listening, but there could be natural contexts where the spectral variation between sounds exceeds the levels that we tested. It is thus possible that spectra occasionally vary enough from sound-to-sound as to necessitate f_0 -based pitch discrimination in real-world conditions. However, this does not appear to be the norm for speech and music sounds, at least for those we had access to here.

Our experiments with natural sounds relied on resynthesis (to render harmonic and inharmonic versions of sounds that were otherwise matched), and this resynthesis might in principle limit the naturalness of the stimuli. The resynthesis is clearly not perfect, but our subjective sense is that the artifacts it induces are modest. Moreover, the extent of the artifacts appears similar for harmonic and inharmonic sounds (the same estimated spectral envelope was used in the resynthesis of harmonic and inharmonic sounds in part to help minimize any such synthesis differences). In particular, speech intelligibility of resynthesized harmonic and inharmonic speech utterances is indistinguishable in quiet (McPherson et al., 2021; Popham et al., 2018). While the synthesis method used in this study was originally designed for speech, the same principles of estimating source vs. filter apply to musical instruments, and the timbre of harmonic instruments remained identifiable after resynthesis. Example vowel and instrument stimuli are available at <https://mcdermottlab.mit.edu/SpectralVariation.html> for readers to judge for themselves.

4.6 The function of f_0 -based pitch

The current results are consistent with a growing body of literature suggesting that in many contexts listeners may use representations of harmonics, rather than the f_0 , for “pitch” discrimination (Faulkner, 1985; McPherson & McDermott, 2018; McPherson & McDermott, 2020; Micheyl et al., 2010; Moore & Glasberg, 1990). The results here provide additional evidence that this phenomenon occurs for natural sounds, in particular speech and musical instruments. The key evidence comes from results with inharmonic stimuli – even though the f_0 of individual inharmonic sounds is ambiguous, the change in individual frequencies from one sound to another is not (at least for modest f_0 differences), and listeners evidently hear this signal as a pitch change and rely on it to make judgments. The results suggest that in many real-world contexts, frequency cues accurately convey f_0 changes, and are used to make judgments about the f_0 . We argue that representations of the harmonics and of the f_0 should both be considered part of pitch perception, as both can be used by listeners to make judgments about f_0 changes depending on the conditions (McPherson & McDermott, 2018; McPherson & McDermott, 2020).

The results here complement evidence that representations of f_0 are used to meet at least two other behavioral challenges: retaining information about sounds over time, and hearing in noise. In other

studies we found harmonic advantages when the sounds being discriminated were either separated by a time delay (McPherson & McDermott, 2020), or superimposed on noise (McPherson et al., 2021). This study adds a third situation in which harmonic advantages are evident – when there are large f_0 differences between sounds passed through a bandpass filter, as is typical for speech and instruments (and many animal vocalizations). By contrast, we found no evidence that f_0 -based pitch enables invariance to differences in spectral shape that occur between natural sounds. The reliance on f_0 representations in settings with extreme spectral variation (e.g. synthetic tones with distinct sets of harmonics) thus appears to be a byproduct of other functions of f_0 -based pitch, including compression for memory, noise robustness, and robust estimation of large f_0 jumps. These findings clarify our understanding of the functional role of pitch.

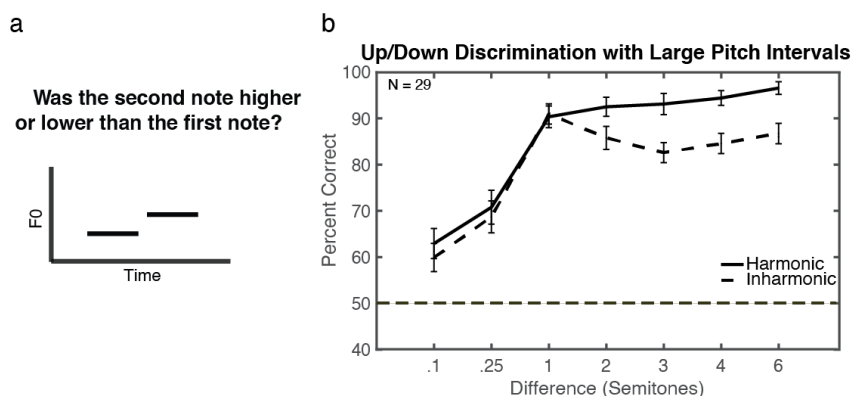
4.7 Future directions

Although we have documented the extent of spectral invariance in pitch perception and clarified the representations that mediate it, we lack a normative understanding of the limits to invariance. It is not obvious why human judgments are not more invariant to differences in the spectral shape of sounds being compared. In particular, listeners do not always adopt strategies that are optimal, at least in the context of certain experimental tasks. For instance, in all experiments we found that discrimination was impaired by spectral differences between sounds, as though listeners are unable to base judgments entirely on the f_0 even when that would maximize task performance. It might be that for natural sounds, changes in f_0 are strongly correlated with changes in the spectral shape (Siedenburg et al., 2021), such that the biases we observed are a signature of an optimal strategy for natural sound discrimination. Recent evidence from infants suggests they may be more robust to spectral differences than adult listeners (Lau et al., 2021), raising the possibility that spectral biases are learned from exposure to natural sounds. Normative models of pitch discrimination (i.e., that instantiate strategies that are optimal under different assumptions) could help to clarify these properties of human pitch perception. Similar questions could be posed about the dependence (or lack thereof) of timbre on pitch (Allen & Oxenham, 2014; Marozeau et al., 2003).

The apparent spectral invariance of pitch representations, in contrast to the limited invariance of pitch judgments, raises the question of where in the brain these effects arise. Proposed neural correlates of f_0 in non-human animals have been reported to be highly invariant to the spectrum in some cases (Bendor & Wang, 2005), and models optimized to estimate f_0 in natural sounds produce relatively invariant representations as well (Saddler et al., 2021). However, auditory cortical responses have also been proposed to underlie pitch judgments (Bizley et al., 2013), and might thus be expected to exhibit spectral biases. And in part because of the coarse scale of human neuroscience methods, pitch-related brain responses in humans have generally not isolated representations of the f_0 (Allen et al., 2017; He & Trainor, 2009; Norman-Haignere et al., 2013; Patterson et al., 2002; Penagos et al., 2004; Tang et al., 2017) leaving the issue open in humans.

Acknowledgements

The authors thank S. Norman-Haignere and L. Demany for helpful discussions and for comments on an earlier version of the manuscript, and M. Saddler for comments on an earlier version of the manuscript. This work was supported by National Institutes of Health (NIH) grants F31DCO18433 and R01DC014739. The funding agency was not otherwise involved in the research, and any opinions, findings, and conclusions or recommendations expressed in this material are those of the authors and do not necessarily reflect the views of the NIH.



Supplementary Figure 1: a. Schematics of task and instructions. b. Results from previous study showing large interval effect (McPherson, M. J., & McDermott, J. H. (2018). Diversity in pitch perception revealed by task dependence. *Nature Human Behavior*, 2, 52-6.). Tones in that study were made inharmonic in the same way that tones were made inharmonic in the current study. The tones in this experiment were filtered with a fixed bandpass filter (Gaussian transfer function on a log-frequency scale, centered at 2,500Hz with a standard deviation of half an octave), applied to ensure that participants could not perform the tasks using changes in the spectral envelope. The tones also had low pass filtered pink noise added to mask distortion products. As in Experiments 1 and 2, with vowels and instrument notes, Harmonic and Inharmonic performance was matched at smaller pitch changes, but performance was impaired for inharmonic tones for larger pitch differences. Error bars show standard error of the mean.

Supplementary Table 1: Vowel Pairs (Vowel 1 vs. Vowel 2 order was randomized in the experiment)

Vowel 1:		Vowel 2:	
IPA		IPA	
/ɛ/	'head'	/i/	'heed'
/i/	'heed'	/u/	'who'd'
/i/	'heed'	/ʊ/	'hood'
/æ/	'had'	/u/	who'd'
/ɔ/	'hod' (like 'cot')	/u/	who'd'
/æ/	'had'	/i/	'heed'
/i/	'heed'	/ʌ/	'hud'
/ɜ~/	'heard'	/i/	'heed'
/ɔ/	'hod' (like 'cot')	/i/	'heed'
/ɑ/	'hawed'	/i/	'heed'

Supplementary Table 2: Instrument Set

1) Alto saxophone
2) Baritone saxophone
3) Bassoon
4) Cello
5) Clarinet
6) Classical guitar
7) Clarinet
8) Cornet
9) Electric guitar
10) English horn
11) Flute
12) French Horn
13) Hammond organ
14) Harpsichord
15) Mandolin
16) Pan flute
17) Piano
18) Pipe organ
19) Soprano saxophone
20) Tenor saxophone
21) Trombone
22) Trumpet
23) Ukulele
24) Viola
25) Violin

Supplementary Table 3: Modulation frequencies and phases for noise amplitude modulation (100% modulation depth)

Low Pass Noise		High Pass Noise	
Frequency (seconds)	Phase	Frequency (seconds)	Phase
.4	π	.4	1
.5	$2.14 (\pi-1)$.5	0
.6	$2.14 (\pi-1)$.6	-.5
.6	$2.64 (\pi-1.5)$.6	-1
.7	1	.7	-1
.7	.5	.7	-1.5
.8	.5	.8	-2
.8	0	.8	-2.5
.9	0	.9	-2.5
.9	-.5	.9	π
1	-.5	1	π
1	-1	1	2.5
1.1	-1	1.1	2.5
1.1	-1.5	1.1	2
1.2	-2.2	1.2	1.5
1.3	π	1.3	1

Supplementary Table 4: Modulation frequencies and phases for noise amplitude modulation for control conditions (100% modulation depth)

Low Pass Noise		High Pass Noise	
Frequency (seconds)	Phase	Frequency (seconds)	Phase
1.1	.5	.5	-2
1.2	0	.6	-2.5
1.3	-.5	1.4	1.3
1.4	-1	1.5	.5
1.6	π	1.6	0

REFERENCES

- Allen, E. J., Burton, P. C., Olman, C. A., & Oxenham, A. J. (2017). Representations of pitch and timbre variation in human auditory cortex. *Journal of Neuroscience*, *37*(5), 1284-1293.
- Allen, E. J., & Oxenham, A. J. (2014). Symmetric interactions and interference between pitch and timbre. *The Journal of the Acoustical Society of America*, *135*(3), 1371-1379.
- Bendor, D., & Wang, X. (2005). The neuronal representation of pitch in primate auditory cortex. *Nature*, *426*, 1161-1165.
- Bianchi, F., Santurette, S., Wendt, D., & Dau, T. (2016). Pitch discrimination in musicians and non-musicians: Effects of harmonic resolvability and processing effort. *Journal of the Association for Research in Otolaryngology*, *17*(1), 69-79.
- Bizley, J. K., Walker, K. M., Nodal, F. R., King, A. J., & Schnupp, J. W. (2013). Auditory cortex represents both pitch judgments and the corresponding acoustic cues. *Current Biology*, *23*(7), 620-625.
- Carruthers, I. M., Laplagne, D. A., Jaegle, A., Briguglio, J. J., Mwilambwe-Tshilobo, L., Natan, R. G., & Geffen, M. N. (2015). Emergence of invariant representation of vocalizations in the auditory cortex. *Journal of Neurophysiology*, *114*(5), 2726-2740.
- Cousineau, M., Carcagno, S., Demany, L., & Pressnitzer, D. (2014). What is a melody? On the relationship between pitch and brightness of timbre. *Frontiers in systems neuroscience*, *7*, 127.
- de Cheveigne, A. (2010). Pitch perception. In C. J. Plack (Ed.), *The Oxford Handbook of Auditory Science: Hearing* (Vol. 3). Oxford University Press.
- DiCarlo, J. J., & Cox, D. D. (2007). Untangling invariant object recognition. *Trends in Cognitive Sciences*, *11*, 333-341.
- Faul, F., Erdfelder, E., Lang, A.-G., & Buchner, A. (2007). G* Power 3: A flexible statistical power analysis program for the social, behavioral, and biomedical sciences. *Behavior research methods*, *39*(2), 175-191.
- Faulkner, A. (1985). Pitch discrimination of harmonic complex signals: Residue pitch or multiple component discriminations? *Journal of the Acoustical Society of America*, *78*, 1993-2004.
- Fletcher, N. H. (1999). The nonlinear physics of musical instruments. *Reports on progress in physics*, *62*(5), 723.
- Fletcher, N. H., & Rossing, T. D. (2010). *The Physics of Musical Instruments*. Springer.
- Glasberg, B. R., & Moore, B. C. J. (1990). Derivation of auditory filter shapes from notched-noise data. *Hearing Research*, *47*, 103-138.
- Goto, M., Hashiguchi, H., Nishimura, T., & Oka, R. (2003). RWC Music Database: Music Genre Database and Musical Instrument Sound Database. The 4th International Conference on Music Information Retrieval (ISMIR 2003),
- Graves, J. E., Micheyl, C., & Oxenham, A. J. (2014). Expectations for melodic contours transcend pitch. *Journal of Experimental Psychology: Human Perception and Performance*, *40*(6), 2338.
- Graves, J. E., Pralus, A., Fornoni, L., Oxenham, A. J., Caclin, A., & Tillmann, B. (2019). Short- and long-term memory for pitch and non-pitch contours: Insights from congenital amusia. *Brain and cognition*, *136*, 103614.
- He, C., & Trainor, L. J. (2009). Finding the pitch of the missing fundamental in infants. *Journal of Neuroscience*, *29*(24), 7718-8822.

- Hillenbrand, J., Getty, L. A., Clark, M. J., & Wheeler, K. (1995). Acoustic characteristics of American English vowels. *Journal of the Acoustical Society of America*, 97(5), 3099-3111.
- JASP. In. (2020). (Version 0.13.1) JASP Team.
- Kawahara, H. (2006). STRAIGHT, exploitation of the other aspect of VOCODER: Perceptually isomorphic decomposition of speech sounds. *Acoustical Science and Technology*, 27(6), 349-353.
- Kawahara, H., & Morise, M. (2011). TANDEM-STRAIGHT: A temporally stable power spectral representation for periodic signals and applications to interference-free spectrum, F0, and aperiodicity estimation. *SADHANA*, 36(5), 713-722.
- Kell, A. J. E., Yamins, D. L. K., Shook, E. N., Norman-Haignere, S., & McDermott, J. H. (2018). A task-optimized neural network replicates human auditory behavior, predicts brain responses, and reveals a cortical processing hierarchy. *Neuron*, 98, 630-644.
- Kishon-Rabin, L., Amir, O., Vexler, Y., & Zaltz, Y. (2001). Pitch discrimination: are professional musicians better than non-musicians? *Journal of Basic Clinical Physiology and Pharmacology*, 12, 125-143.
- Kleiner, M., Brainard, D., & Pelli, D. (2007). What's new in Psychtoolbox-3?
- Lau, B. K., Oxenham, A. J., & Werner, L. A. (2021). Infant Pitch and Timbre Discrimination in the Presence of Variation in the Other Dimension. *Journal of the Association for Research in Otolaryngology*, 1-10.
- Licklider, J. C. R. (1954). "Periodicity" Pitch and "Place" Pitch. *The Journal of the Acoustical Society of America*, 26(5), 945-945.
- Liu, S. T., Montes-Lourido, P., Wang, X., & Sadagopan, S. (2019). Optimal features for auditory categorization. *Nature Communications*, 10, 1302.
- Marozeau, J., de Cheveigné, A., McAdams, S., & Winsberg, S. (2003). The dependency of timbre on fundamental frequency. *The Journal of the Acoustical Society of America*, 114(5), 2946-2957.
- McDermott, J. H., Ellis, D. P. W., & Kawahara, H. (2012). Inharmonic speech: A tool for the study of speech perception and separation. Proceedings of SAPA-SCALE, Portland, OR.
- McDermott, J. H., Keebler, M. V., Micheyl, C., & Oxenham, A. J. (2010). Musical intervals and relative pitch: Frequency resolution, not interval resolution, is special. *The Journal of the Acoustical Society of America*, 128(4), 1943-1951.
- McDermott, J. H., Lehr, A. J., & Oxenham, A. J. (2008). Is relative pitch specific to pitch? *Psychological Science*, 19(12), 1263-1271.
- McPherson, M. J., Dolan, S. E., Durango, A., Ossandon, T., Valdés, J., Undurraga, E. A., Jacoby, N., Godoy, R. A., & McDermott, J. H. (2020). Perceptual fusion of musical notes by native Amazonians suggests universal representations of musical intervals. *Nature Communications*, 11, 2786.
- McPherson, M. J., Grace, R. C., & McDermott, J. H. (2021). Harmonicity aids hearing in noise. *Attention, Perception, and Psychophysics* (In Press).
- McPherson, M. J., & McDermott, J. H. (2018). Diversity in pitch perception revealed by task dependence. *Nature Human Behavior*, 2, 52-66.
- McPherson, M. J., & McDermott, J. H. (2020). Time-dependent discrimination advantages for harmonic sounds suggest efficient coding for memory. *Proceedings of the National Academy of Sciences*, 117(50), 32169-32180.

- McWalter, R., & McDermott, J. H. (2019). Illusory sound texture reveals multi-second statistical completion in auditory scene analysis. *Nature Communications*, *10*, 5096.
- Melara, R. D., & Marks, L. E. (1990). Interaction among auditory dimensions: Timbre, pitch, and loudness. *Perception & Psychophysics*, *48*(2), 169-178.
- Micheyl, C., Delhommeau, K., Perrot, X., & Oxenham, A. J. (2006). Influence of musical and psychoacoustical training on pitch discrimination. *Hearing Research*, *219*, 36-47.
- Micheyl, C., Divis, K., Wroblewski, D. M., & Oxenham, A. J. (2010). Does fundamental-frequency discrimination measure virtual pitch discrimination? *Journal of the Acoustical Society of America*, *128*(4), 1930-1942.
- Micheyl, C., & Oxenham, A. J. (2004). Sequential F0 comparisons between resolved and unresolved harmonics: No evidence for translation noise between two pitch mechanisms. *The Journal of the Acoustical Society of America*, *116*.
- Moore, B. C., Glasberg, B. R., & Proctor, G. M. (1992). Accuracy of pitch matching for pure tones and for complex tones with overlapping or nonoverlapping harmonics. *The Journal of the Acoustical Society of America*, *91*(6), 3443-3450.
- Moore, B. C. J., & Glasberg, B. R. (1990). Frequency discrimination of complex tones with overlapping and non-overlapping harmonics. *Journal of the Acoustical Society of America*, *87*, 2163-2177.
- Norman-Haignere, S., Kanwisher, N., & McDermott, J. H. (2013). Cortical pitch regions in humans respond primarily to resolved harmonics and are located in specific tonotopic regions of anterior auditory cortex. *Journal of Neuroscience*, *33*(50), 19451-19469.
- Norman-Haignere, S., & McDermott, J. H. (2016). Distortion products in auditory fMRI research: Measurements and solutions. *NeuroImage*, *129*, 401-413.
- Patterson, R. D., Uppenkamp, S., Johnsrude, I. S., & Griffiths, T. D. (2002). The processing of temporal pitch and melody information in auditory cortex. *Neuron*, *36*(4), 767-776.
- Penagos, H., Melcher, J. R., & Oxenham, A. J. (2004). A neural representation of pitch salience in nonprimary human auditory cortex revealed with functional magnetic resonance imaging. *Journal of Neuroscience*, *24*(30), 6810-6815.
- Plack, C. J., Oxenham, A. J., Popper, A. J., & Fay, R. R. (Eds.). (2005). *Pitch: Neural Coding and Perception*. Springer.
- Popham, S., Boebinger, D., Ellis, D. P., Kawahara, H., & McDermott, J. H. (2018). Inharmonic speech reveals the role of harmonicity in the cocktail party problem. *Nature Communications*, *9*(1), 2122.
- Pressnitzer, D., & Patterson, R. D. (2001). Distortion products and the perceived pitch of harmonic complex tones. In D. J. Breebaart (Ed.), *Physiological and Psychophysical Bases of Auditory Function* (pp. 97-104). Shaker Publishing.
- Rouder, J. N., Speckman, P. L., Sun, D., Morey, R. D., & Iverson, G. (2009). Bayesian t tests for accepting and rejecting the null hypothesis. *Psychonomic Bulletin & Review*, *16*(2), 225-237.
- Russo, F. A., & Thompson, W. F. (2005). An interval size illusion: The influence of timbre on the perceived size of melodic intervals. *Perception & Psychophysics*, *67*(4), 559-568.
- Saddler, M. R., Gonzalez, R., & McDermott, J. H. (2021). Deep neural network models reveal interplay of peripheral coding and stimulus statistics in pitch perception. *Nature Communications*, *12*(7278).
- Semal, C., & Demany, L. (1991). Dissociation of pitch from timbre in auditory short-term memory. *The Journal of the Acoustical Society of America*, *89*(5), 2404-2410.

- Sharpee, T. O., Atencio, C. A., & Schreiner, C. E. (2011). Hierarchical representations in the auditory cortex. *Current Opinion in Neurobiology*, 21(5), 761-767.
- Siedenburg, K. (2018). Timbral Shepard-illusion reveals ambiguity and context sensitivity of brightness perception. *The Journal of the Acoustical Society of America*, 143(2), EL93-EL98.
- Siedenburg, K., Jacobsen, S., & Reuter, C. (2021). Spectral envelope position and shape in sustained musical instrument sounds. *The Journal of the Acoustical Society of America*, 149(6), 3715-3726.
- Singh, P. G., & Hirsh, I. J. (1992). Influence of spectral locus and F 0 changes on the pitch and timbre of complex tones. *The Journal of the Acoustical Society of America*, 92(5), 2650-2661.
- Slaney, M. (1998). Auditory toolbox. *Interval Research Corporation, Tech. Rep*, 10.
- Spiegel, M. F., & Watson, C. S. (1984). Performance on frequency-discrimination tasks by musicians and non-musicians. *Journal of the Acoustical Society of America*, 76, 1690-1695.
- Stevens, K. N. (2000). *Acoustic Phonetics*. MIT Press.
- Tang, C., Hamilton, L. S., & Chang, E. F. (2017). Intonational speech prosody encoding in the human auditory cortex. *Science*, 357(6353), 797-801.
- Traer, J., Norman-Haignere, S. V., & McDermott, J. H. (2021). Causal inference in environmental sound recognition. *Cognition*, 214, 104627.
- Vurma, A., Raju, M., & Kuuda, A. (2011). Does timbre affect pitch?: Estimations by musicians and non-musicians. *Psychology of Music*, 39(3), 291-306.
- Warrier, C. M., & Zatorre, R. J. (2002). Influence of tonal context and timbral variation on perception of pitch. *Perception & Psychophysics*, 64(2), 198-207.
- Woods, K. J. P., & McDermott, J. (2018). Schema learning for the cocktail party problem. *Proceedings of the National Academy of Sciences*, 115, E3313-E3322.
- Woods, K. J. P., Siegel, M. H., Traer, J., & McDermott, J. H. (2017). Headphone screening to facilitate web-based auditory experiments. *Attention, Perception, and Psychophysics*, 79, 2064-2072.
- Zarate, J. M., Ritson, C. R., & Poeppel, D. (2013). The effect of instrumental timbre on interval discrimination. *PLoS One*, 8(9), e75410.