

Elsevier Editorial System(tm) for Journal of  
Phonetics  
Manuscript Draft

Manuscript Number:

Title: Hidden dimensions in phonetic spaces: temporal resolution in Cantonese tone perception

Article Type: Research Article

Keywords: tone perception; tone recognition; temporal integration; temporal resolution; Cantonese; interrupted speech

Corresponding Author: Ms. Kristine Yu,

Corresponding Author's Institution:

First Author: Kristine Yu

Order of Authors: Kristine Yu

Abstract: The role of temporal resolution in speech perception (e.g. whether tones are parameterized with fundamental frequency sampled every 10 ms, or just twice in the syllable) is little studied, and the temporal resolution relevant for tonal perception is still an open question. The choice of temporal resolution matters because how we understand the recognition, dispersion, and learning of phonetic categories is entirely predicated on what dimensions we use to define the phonetic space that they lie in. Here, we present a tonal perception experiment in Cantonese where we used interrupted speech in trisyllabic stimuli to study the effect of temporal resolution on human tonal identification. We also performed acoustic classification of the stimuli with support vector machines. Our results show that just a few samples per syllable are enough for humans and machines to classify Cantonese tones with reasonable accuracy, without much difference in performance from having the full speech signal available. However, the identification of a subset of tones does fail under coarse sampling, and where in the syllable samples are taken is critical: sampling at the syllable onset is critical for identifying rising tones, due to peak delay.



UNIVERSITY OF MASSACHUSETTS  
AMHERST

Integrative Learning Center  
650 North Pleasant Street  
Amherst, MA 01003-1100

Department of Linguistics

voice: 413.545.0885  
fax: 413.545.2792  
[www.umass.edu/linguistics](http://www.umass.edu/linguistics)

June 13, 2016

Dear Editors of Journal of Phonetics,

I am submitting for review an article entitled “Hidden dimensions in phonetic spaces: temporal resolution in Cantonese tone perception”, with the following abstract:

*The role of temporal resolution in speech perception (e.g. whether tones are parameterized with fundamental frequency sampled every 10 ms, or just twice in the syllable) is little studied, and the temporal resolution relevant for tonal perception is still an open question. The choice of temporal resolution matters because how we understand the recognition, dispersion, and learning of phonetic categories is entirely predicated on what dimensions we use to define the phonetic space that they lie in. Here, we present a tonal perception experiment in Cantonese where we used interrupted speech in trisyllabic stimuli to study the effect of temporal resolution on human tonal identification. We also performed acoustic classification of the stimuli with support vector machines. Our results show that just a few samples per syllable are enough for humans and machines to classify Cantonese tones with reasonable accuracy, without much difference in performance from having the full speech signal available. However, the identification of a subset of tones does fail under coarse sampling, and where in the syllable samples are taken is critical: sampling at the syllable onset is critical for identifying rising tones, due to peak delay.*

This manuscript is not and will not be submitted elsewhere prior to an editorial decision.

I look forward to the review process. Thank you for your consideration.

Sincerely,

Kristine Yu  
Assistant Professor  
Department of Linguistics  
University of Massachusetts Amherst

1      Hidden dimensions in phonetic spaces: temporal resolution in Cantonese tone perception  
 2  
 3  
 4  
 5  
 6  
 7  
 8  
 9

---

10     Anonymous  
 11  
 12  
 13  
 14  
 15  
 16  
 17  
 18  
 19  
 20  
 21  
 22  
 23  
 24  
 25

## Abstract

The role of temporal resolution in speech perception (e.g. whether tones are parameterized with fundamental frequency sampled every 10 ms, or just twice in the syllable) is little studied, and the temporal resolution relevant for tonal perception is still an open question. The choice of temporal resolution matters because how we understand the recognition, dispersion, and learning of phonetic categories is entirely predicated on what dimensions we use to define the phonetic space that they lie in. Here, we present a tonal perception experiment in Cantonese where we used interrupted speech in trisyllabic stimuli to study the effect of temporal resolution on human tonal identification. We also performed acoustic classification of the stimuli with support vector machines. Our results show that just a few samples per syllable are enough for humans and machines to classify Cantonese tones with reasonable accuracy, without much difference in performance from having the full speech signal available. However, the identification of a subset of tones does fail under coarse sampling, and where in the syllable samples are taken is critical: sampling at the syllable onset is critical for identifying rising tones, due to peak delay.

**Keywords:** tone perception, tone recognition, temporal integration, temporal resolution, Cantonese, interrupted speech

## 1. Introduction

A central question of phonetics is understanding how phonetic concepts are defined (Ladefoged, 1980). A phonetic concept picks out a subset of some space defined over phonetic dimensions. If we define phonetic dimensions to be any property relevant to the human determination of speech sounds, then they include contributions from visual (e.g. McGurk & MacDonald (1976)) and top-down contextual (e.g. Ganong III (1980)) properties. Narrowly speaking, though, phonetic dimensions are often thought as acoustic, articulatory, or auditory, and this subset of dimensions is what we focus on here. To understand how a phonetic concept—e.g. the high tone ˥ in Cantonese, or the [ae] vowel in English—is defined, a critical question to ask is: what are the dimensions of the phonetic space in which the concept is defined? One way to operationalize this question further is to ask: what are the primary<sup>1</sup> measures or parameters that distinguish contrasting phonetic concepts? For example, one familiar set of acoustic dimensions for defining English vowels is the first and second formants (measured at steady state) (Peterson & Barney, 1952).

The high-level contribution of this paper is to draw attention to a hidden source of dimensions in phonetic spaces: temporal resolution. The speech signal unfolds in time, so any parameter, e.g. the first formant (F1) or fundamental frequency (f0) actually expands into a *family* of parameters. Each parameter is

sampled over time, e.g. if F1 is parameterized to be measured at the onset and offset of a vowel, then the F1 family contributes 2 dimensions to the defined vowel space; if f0 is parameterized to be measured at 8 timepoints for a tone, then f0 contributes 8 dimensions to the defined tonal space.<sup>2</sup>

Why does it matter if we define an 8-dimensional tonal space with 8 f0 samples over the syllable (Khouw & Ciocca, 2007) vs. only a 2-dimensional one, with f0 samples measured at the onset and offset of voicing (Barry & Blamey, 2004)? It matters because how we understand the recognition, dispersion, and learning of tonal concepts is entirely predicated on what dimensions we use to define the phonetic space that they lie in. For example, Kuang (2013) showed that once phonation parameters in addition to f0 parameters are included as dimensions in a tonal space, we can understand how listeners can possibly discriminate between the five level tones of Black Miao. Choosing the temporal resolution for f0 can have as much import as deciding whether or not non-f0 properties should be included as parameters in a tonal space: Alexander (2010) found that tones were not always well-dispersed within an inventory across a range of languages if the tonal space was defined using a single f0 point; however tones were well-dispersed if they were defined in a 2-D tonal space over f0 measured at the onset and offset of the syllable. Choosing how to parameterize f0 properties also matters; Gauthier et al. (2007) found that neural networks learned to classify Mandarin tones with higher accuracy when the tonal space was defined over the f0 velocity timecourse than the f0

<sup>1</sup>Even setting aside non-acoustic properties, there are already an infinite number of parameters one could extract from the speech signal, and the primary cues that listeners use to identify sound categories may be large (Lisker, 1978) and flexible (Whalen & Xu, 1992; Liu & Samuel, 2004; Sumner & Samuel, 2009; Clarke & Garrett, 2004). Thus, what one means by “primary” must also be operationalized, but no matter what the criteria, we must have some way of limiting the number of parameters in play for scientific interpretability.

<sup>2</sup>The shape of the f0 curve may also be parameterized in terms of a family of functions (Hirst & Espesser, 1993; Taylor, 2000; Andruski & Costello, 2004; Kochanski et al., 2005; Hermes, 2006; Prom-on et al., 2009; Shih & Lu, 2015; Li & Chen, 2016), such as quadratic polynomials, but the issue of temporal resolution remains. For instance, the more finely one wishes to capture the detailed shape of the contour, the higher the degree of polynomial needed.

1 timecourse.

2 This paper uses Cantonese tone perception as a case study  
3 to examine the contribution of fine temporal detail to the com-  
4 plexity of phonetic spaces. We use perceptual and acoustic ev-  
5 idence to address how humans and machines respond in tone  
6 identification as the temporal resolution in the speech signal is  
7 systematically lowered.

### 8   9 *1.1. Temporal resolution for tonal concepts*

10 The degree of fineness of temporal resolution in the speech  
11 signal relevant for human cognition is still an open question.  
12 It has long been assumed in linguistics that fine temporal res-  
13 olution of the speech signal is not necessary for the parame-  
14 terization of tones in tone languages, and discussions of tem-  
15 poral resolution have largely been confined to the automatic  
16 tonal recognition literature. Chao, who introduced the iconic  
17 tone letters (Chao, 1930) used in the International Phonetic Al-  
18 phabet for representing linguistic tone, wrote: “the exact shape  
19 of the time-pitch curve, so far as I have observed, has never  
20 been a necessary distinctive feature, given the starting and end-  
21 ing points, or the turning point, if any, on the five-point scale”  
22 (Chao, 1968, 25), and tone letters are understood to have up to  
23 3 samples, e.g. ↘.

24 Additionally, Laniran (1992) argued for two targets per tone  
25 in Yoruba, and Barry & Blamey (2004) argued for 2-D acous-  
26 tic Cantonese tonal spaces defined over onset and offset f0 val-  
27 ues based on perceptual dimensions hypothesized from mul-  
28 tidimensional scaling analyses of cross-linguistic tonal per-  
29 ception (Gandour & Harshman, 1978; Gandour, 1981, 1983).  
30 Morén & Zsiga (2006) and Zsiga & Nitisoroj (2007) argued for  
31 one target per tone in the representation of Thai tones in con-  
32 nected speech, with falling tones associated with a peak at syl-  
33 lable midpoint, high tones with a peak at syllable offset, rising  
34 tones with an f0 minimum at syllable midpoint, and low tones  
35 with a low pitch target at syllable offset.

36 In contrast, the computational literature has sometimes pre-  
37 sumed that fine sampling is valuable. In a study of unsuper-  
38 vised learning of Mandarin tones, Gauthier et al. (2007) ex-  
39 tracted 30 samples of f0 or 28 samples of f0 velocity per syllable  
40 (a sampling rate on the order of 1 sample every 10ms), and  
41 a number of automatic tonal recognizers parameterize the f0  
42 curve by sampling f0 every 10ms, e.g. Zhang & Hirose (2004);  
43 Pisarn & Theeramunkong (2006); Prukkanon et al. (2016). But  
44 in support of the classic linguistic intuition of sparse tempo-  
45 ral resolution in representing tones, Tian et al. (2004)’s auto-  
46 matic tonal recognition study of Mandarin previously showed  
47 that sparse temporal resolution, with 4 samples/tone, can out-  
48 perform fine-grained sampling with 1 sample/10 ms, concluding  
49 that “detailed information is useless for tone discrimina-  
50 tion” (Tian et al., 2004, I-107).

51 Other Mandarin and Cantonese tonal recognizers have used  
52 a simple time warping-like (time normalization) coarse sam-  
53 pling scheme of 3-5 f0 averages or frame values over uniformly  
54 divided subsegments of (part of) the syllable (Peng & Wang,  
55 2005; Qian et al., 2007; Wang & Levow, 2008; Zhou et al.,  
56 2008). For the synthesis of natural-sounding Cantonese tones,  
57 Li & Lee (2007) and Li & Lee (2008) argued that one or two  
58

59 linear movements per tone sufficed. In sum, the computational  
60 literature on tone recognition has not settled on the fineness of  
61 sampling resolution to use in tonal feature extraction from the  
62 speech signal, although the common use of low temporal reso-  
63 lution in feature extraction for tonal recognition is striking com-  
64 pared to the dominance of fine-grained sampling in the rest of  
65 the speech recognition literature.

66 But recent work on tonal production, perception, and pro-  
67 cessing has raised questions about the assumption of the suf-  
68 ficiency of coarse sampling for tones. Barnes et al. (2012)  
69 argued that details of contour shape such as convexity and  
70 concavity define English intonational pitch accent contrasts.  
71 Remijsen (2013); Remijsen & Ayoker (2014) discovered con-  
72 trastive early/late falls in Dinka and Shilluk, and DiCanio et al.  
73 (2014) found contrastive early/late rises in Yoloxóchitl Mixtec.  
74 In Dinka and Shilluk, a difference of a mere 40-64 ms in the  
75 onset of an f0 fall signals the lexical contrast between an early  
76 fall and a late fall. Remijsen (2013) also provided evidence that  
77 Dinka speakers had no trouble discriminating between the early  
78 and late falls in perception. Moreover, Chandrasekaran et al.  
79 (2007) and Krishnan et al. (2009) found in electroencephalo-  
80 graphic studies that Chinese speakers showed a language ex-  
81 perience advantage relative to English speakers in processing  
82 rising f0 contours only when their curvature matched the con-  
83 cave shape of the rising tone in Chinese, and not when they  
84 were approximations of one or two linear segments or convex.  
85 In sum, there is a growing body of evidence supporting the need  
86 for fine temporal resolution to capture details of the f0 contour  
87 shape to distinguish tonal contrasts.

### 88   9 *1.2. Goals of current study*

89 While the evidence presented in the previous section indi-  
90 rectly bears on the fineness of temporal resolution in tonal per-  
91 ception, there is little work that explicitly tests the effect of  
92 temporal resolution on tone perception. Our study does just  
93 that. Perceptual studies that present listeners with tones resyn-  
94 thesized with particular restrictive parameterizations, e.g. as  
95 polynomials of a certain degree, can provide at least a pre-  
96 liminary indication of whether the restrictive parameterization  
97 chosen could be a close approximation to that in human percep-  
98 tion by collecting similarity judgments between the resyntheses  
99 and original stimuli (Hermes, 2006; Li & Lee, 2007). Here, we  
100 backed off from imposing a hypothesized restrictive parameter-  
101 ization on the listener and focused on the more general issue of  
102 assessing the effect of reducing the temporal resolution of the  
103 speech signal on tonal perception. We manipulated which time-  
104 points in the stimuli the listeners had the opportunity to hear by  
105 intermittently deleting the recorded speech signal and replac-  
106 ing it with white noise, in the tradition of phoneme restoration  
107 (Miller & Licklider, 1950; Warren, 1970; Bashford et al., 1992;  
108 Samuel, 1996) and “silent center” (Strange et al., 1983) percep-  
109 tual experiments.

110 Gottfried & Suiter (1997) and Lee et al. (2008, 2009); Lee  
111 (2009) previously performed “silent center” Mandarin tonal  
112 perception experiments, which could be construed as manip-  
113 ulations of temporal resolution in the speech signal. Lee et al.  
114 (2008, 2009); Lee (2009) built on Gottfried & Suiter (1997)’s

1 small-scale study and included comparisons of tonal identification accuracy for listeners under time pressure between intact  
2 tones, “silent center” (Strange et al., 1983) tones with only initial and final regions available and the speech material in between silenced (“2 samples” over the tone; the first 6 and final 8 pitch periods), and tones with only the initial region available (“1 sample” over the tone; the first 6 pitch periods). Key results from these studies are that: (a) tonal identification accuracy did decrease as a function of the amount of input available to the listener, but remained high and well above chance (25%)—mostly 80%–95% accuracy—regardless of whether the stimuli were single speaker or multispeaker, whether presented in isolation or with preceding context, and whether the preceding context was cross-spliced in from another recording or not; (b) providing preceding context significantly facilitated tonal identification for the onset-only and silent-center conditions, relative to providing no context; (c) reducing the amount of speech signal available to the listener affected different tones differently.

Our study builds on this work. First, we chose to perform the study in Cantonese rather than Mandarin. Tonal identification in Cantonese presents a more challenging tonal identification task than in Mandarin, so we can avoid the ceiling effects that occurred with Mandarin. While Mandarin has four tones which all have very different f0 contour shapes, the tonal inventory of Cantonese includes three level tones (high level Tone 1, T55, ↑; mid level Tone 3, 33, +; low level Tone 6, T22, ↓), two rising tones (high rising Tone 2, T25, ↗; low rising Tone 5, T23, ↘), and a falling tone (Tone 4, T21, ↘), cf. Figure 2 (Matthews & Yip, 1994).<sup>3</sup>

Second, we make a more explicit connection between what acoustic information is available in the signal and how listeners perceived tones as the signal is degraded. To do this, we accompanied the human perception experiment with a machine classification study. To make the human and machine results comparable, we designed our study to try to simulate the conditions of machine classification in the perception task for the human listeners. We provided an experimental context for tonal identification limited in a way to be similar to characteristics of feature extraction in automatic tonal recognition. We used tritone stimuli from connected speech, as most recent automatic tonal recognizers use acoustic feature extraction from a temporal window extending beyond a single tone to its neighbors (Zhang & Hirose, 2000; Levow, 2005; Qian et al., 2007), and we used stimuli from multiple speakers like in the speaker-independent tonal recognition tasks in Peng & Wang (2005); Qian et al. (2007). We also resynthesized the syllable durations of the tritones to be fixed at their grand average to simulate the commonly employed preprocessing step of time normalization to the syllable. Our experimental manipulation of temporal res-

olution in the signal used interrupting noise to create a 5-step gradient of sampling resolution and make uniformly distributed “samples” or windows from the speech signal available to the listener—a very simple treatment designed to simulate the common uniformly sampled vector time series feature extraction procedure in automatic tonal recognition.

Finally, our study connects how often the speech signal is sampled to *where* the speech signal is sampled. Khouw & Ciocca (2007) found that f0 information at the syllable offset was the most critical in acoustic and perceptual discrimination of Cantonese tones in isolation. To follow up on this, while Lee et al. (2008, 2009)’s had a single 1-sample condition sampling from the target syllable onset, we compare sampling from the syllable onset, midpoint, and offset in our machine classification experiments. Also, while Lee et al. (2008, 2009) only preceded the target syllable with other speech material, we include a syllable *following* the target syllable, in addition to one preceding the syllable, as in Gottfried & Suiter (1997). This following syllable provides a buffer for f0 information in the target syllable offset that may be shifted or spread onto the following syllable.

Based on the past work discussed in this section, we had three general hypotheses for our study:

1. *Sufficiency of coarse resolution*: like Lee et al. (2008, 2009) found for Mandarin, tonal identification accuracy by humans and machine is well above chance for all temporal resolution conditions, with little detriment to identification accuracy overall as the sampling becomes coarser.
2. *Tone-specific disadvantage with coarser resolution*: the brunt of the deterioration in tonal identification as resolution drops comes from confusability involving the two rises, T25 and T23, which participate in subtle contrasts in f0 contour shape similar to that of the contrastive falls/rises found in Remijzen (2013); Remijzen & Ayoker (2014); DiCanio et al. (2014).
3. *Informativity of syllable offset*: based on Khouw & Ciocca (2007), the availability of acoustic information from syllable offsets facilitates perception more than information from syllable onsets or midpoints.

In the rest of this paper, we describe the speech materials used in the perception experiment and procedures for the experiment and analysis (§2), present results from the perception experiment and machine classification task (§3), discuss these results (§4), and conclude in §5.

## 2. Materials and methods

### 2.1. Recordings

The stimuli were recorded by ten native Cantonese speakers, five of whose recordings were further processed for the rest of the study: these three males and two females were chosen to span a wide pitch range (see Appendix A), to provide a representative instance of the challenge of a multispeaker task. Four of the speakers were born and raised in Hong Kong and recorded in the phonetics lab sound-attenuated booth at

<sup>3</sup>Descriptions vary slightly in the exact 5-value integers assigned to the tones, but the exact integers used here are not of importance since we use these designations purely as mnemonic names. Some descriptions also distinguish these tones from the shorter entering tones (high, mid, and low level) which occur in syllables with unreleased stop codas. Throughout the paper, we use 5-valued integer designations as mnemonic names for the tonal categories, e.g. T55 for ↑.

1 the City University of Hong Kong. One was born and raised  
2 in Macau and recorded in the phonetics lab sound-attenuated  
3 booth at University of California, Los Angeles. They were  
4 recruited from the local university student population and re-  
5 ceived cash compensation. All speakers were recorded using a  
6 Shure SM10A-CN headworn mic. For the speaker at UCLA,  
7 the signal was run through an XAudioBox pre-amplifier and A-  
8 D device to a computer at 22,050 Hz/16 bits with PCQuirerX  
9 (Scicon R&D, Inc.). The speakers in Hong Kong were recorded  
10 at 44.1kHz/16 bits with a TASCAM HD-R1 digital recorder.

11 The stimuli were created from the tritone  
12 ⟨wai+, {wai˥, ˥˧, ˧˥, ˧˩, ˩˧, ˩˩}, mat+⟩ (wai<sup>33</sup> wai mat<sup>3</sup>) ex-  
13 tracted from sentences of the form: lei<sup>25/35</sup> yiu<sup>33</sup> wai<sup>33</sup> wai  
14 mat<sup>3</sup> deng/geng<sup>33</sup> ‘you want Wai-Wai to clean the lamp/mirror’  
15 with the target, the second /wai/, ranging over all six Cantonese  
16 tones. They were part of a larger study on contextual tonal  
17 variability. The lexical meanings of the orthographic characters  
18 we associated with tones T55, T25, T33, T21, T23, and T22  
19 were, respectively, ‘power’, ‘appoint’, ‘fear’, ‘surround’,  
20 ‘great’, and ‘stomach’, and speakers were asked to treat /wai  
21 wai/ as a (nonce) proper name. The orthographic characters  
22 were chosen to be the most familiar ones for each tone by  
23 a native speaker. Each speaker actually recorded 5 fluent  
24 repetitions of sentences containing all 36 bitone combinations  
25 over /wai wai/ (with the sentences not used as stimuli for  
26 the perception experiment serving as fillers), from which we  
27 chose the last three repetitions of each Tone 33-Tone X-Tone 3  
28 tritone for the stimuli set for a total of 90 tritones, 18 from each  
29 speaker, 3 distinct repetitions per speaker per tritone.<sup>4</sup> We held  
30 the initial and final syllable of the tritone constant as mid-level  
31 tones so that the listener’s task was limited to the identification  
32 of a single tone, since a pilot experiment allowing variation  
33 in more than one tone was very confusing to participants. A  
34 Cantonese native speaker trained in linguistics and phonetics  
35 checked that none of the speakers had tonal mergers and that  
36 the speakers uttered the tones correctly. No speakers produced  
37 Tone T55 with a 53 high fall contour, a variant more common  
38 in the past.

## 41 42 2.2. Resynthesis 43

44 All stimuli were resampled to 22kHz; tritones were ex-  
45 tracted using a rectangular window, and RMS amplitude was  
46 rescaled to 75 dB (relative to the auditory threshold) in Praat  
47 (Boersma & Weenink, 2010). All syllables were resynthesized  
48 using PSOLA implemented in Praat to be have a target dura-  
49 tion of 241 ms, the grand mean of the syllable durations, for a  
50 total duration of 740 ms for the tritone, to simulate time normal-  
51 ization to the syllable.<sup>5</sup> The manipulated condition, TEMPORAL  
52 RESOLUTION, was varied from the intact signal, to 7, 5, 3, and  
53 2 uniformly spaced samples (time-slices or windows) of 30.41  
54 ms each per syllable. The sample duration was well below the  
55

56  
57 <sup>4</sup>In three cases, we chose another repetition than those listed above due to  
58 sound quality of the recording.

59 <sup>5</sup>The PSOLA algorithm resynthesis added about 18 ms over the target dura-  
60 tion over the course of the tritone.

minimum 130 ms duration Greenberg & Zee (1977) found nec-  
essary for perception of a nonzero f0 velocity, “contouricity”,  
in speech, and also on the same order of magnitude as the stand-  
ard frame size in automated short-term analysis f0 detection  
(Hess, 1983, 343).

The temporal resolution manipulation involved intermit-  
tently deleting the recorded speech signal and replacing it with  
white noise low-pass-filtered at 5000 Hz that was 10dB higher  
than the average signal amplitude, cf. Figure 1. Similar stim-  
uli manipulations are used in phonemic restoration studies, in  
which listeners perceive segmental speech sounds to be present  
in the presence of noise even if they are not (Miller & Licklider,  
1950; Warren, 1970; Bashford et al., 1992; Samuel, 1996). We  
alternated the speech signal with louder noise rather than si-  
lence because the intelligibility of the speech is well-known to  
be poor when alternated with silent gaps; however, continuity  
of the speech percept can be maintained when the speech sig-  
nal is alternated with a louder sound that is a potential masker  
of the fainter speech signal. This phenomenon is in fact the  
basis of phonemic restoration. Broadband noise has typically  
been used in segmental phoneme restoration experiments, and  
we chose to use white noise low-pass-filtered at 5000 Hz in  
particular because it has also been used in studying the conti-  
nuity of tones through interrupting noise (Ciocca & Bregman,  
1987). Additionally, we chose white noise to avoid providing  
any information that the listener might use in perceiving the in-  
terrupted speech, since Bashford et al. (1996) showed a boost  
in the intelligibility of speech interrupted by speech-modulated  
noise rather than white noise.

The noise was generated using the MLP Matlab tool-  
box (Grassi & Soranzo, 2009). Since the sample durations were  
fixed, the noise duration varied for different RESOLUTION con-  
ditions, but was fixed within a condition, ranging from 90ms  
to 4ms from the 2- to 7- sample condition, respectively, as  
shown in Figure 1. The duration of noise intervals was mea-  
sured in absolute time rather than the number of glottal pulses,  
cf. Lee et al. (2009), to simulate the constant frameshift used in  
feature extraction in automatic tone recognizers. The noise in-  
tervals included raised-cosine onset and offset ramps that were  
10% of the duration of the noise interval to reduce audible spec-  
tral splatter (Hant & Alwan, 2003); the duration of the ramps  
was chosen to be relative to the duration of the noise interval  
since the noise interval duration varied between sampling res-  
olution conditions. Half-duration noise intervals were used at  
the onset and offset of the tritone, with extra noise padding at  
the offset if needed to replace the entirety of the duration of  
the intact speech signal. Due to a programming error not de-  
tected until after the participants were tested, the last noise in-  
terval for the 2- and 3-sample stimuli was of full rather than  
half duration. For these two conditions, the final noise interval  
thus extended the stimulus duration beyond that of the stimuli  
for the other conditions. However, the same information from  
the speech signal was available to the listeners that would have  
been present without the added noise: the extended noise at the  
stimulus offset did not replace any speech information.

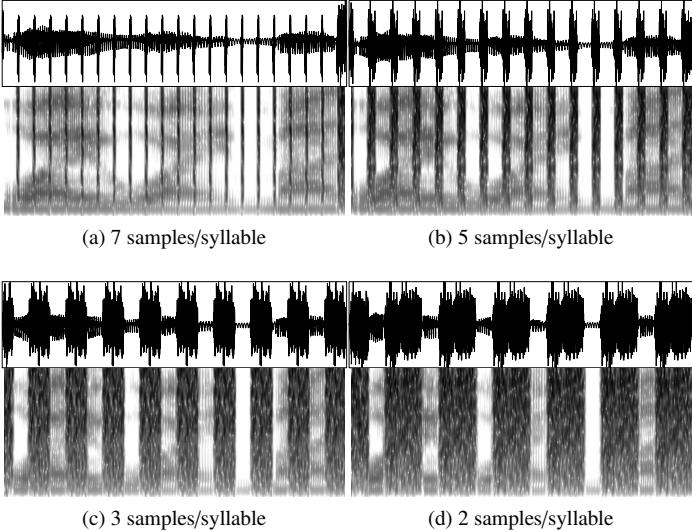


Figure 1: Waveforms and spectrograms of a Cantonese tritone Tone T33 - Tone T21 - Tone 33 stimulus under different temporal resolution conditions from intact, to 7, 5, 3, and 2 samples/syllable.

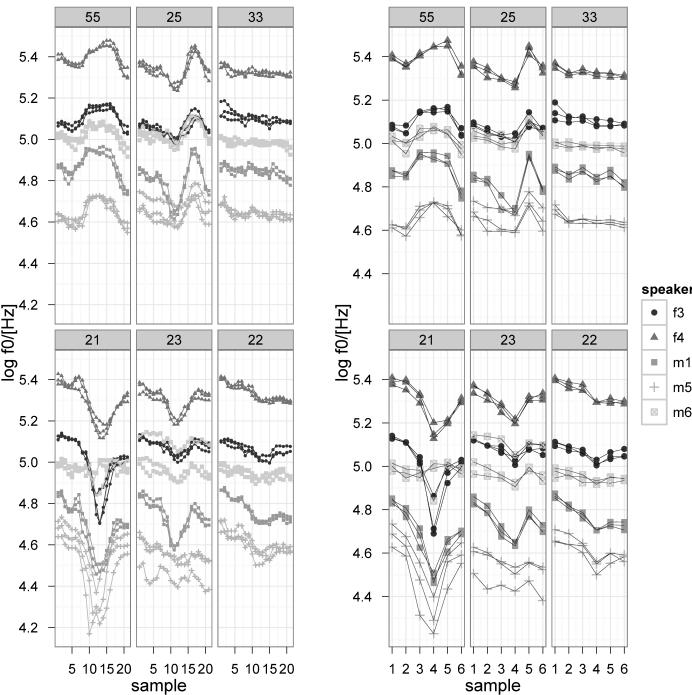


Figure 2: Fundamental frequency (f0) contours extracted with RAPT using speaker-specific pitch floors and ceilings, showing the parameterization of f0 contours for the 7-sample and 2-sample conditions for computational modeling. Linear interpolation was used for replacing missing values and smoothing. (left panel) log-transformed f0 extracted with 10ms frameshift and averaged over each of the 21 samples in the 7 samples/syllable condition. (right panel) log-transformed f0 values averaged over each of 6 samples in the 2 samples/syllable condition.

### 2.3. Acoustic feature extraction

Fundamental frequency (f0) timecourses were extracted for stimuli description and exploratory machine classification, cf. §2.6.3. The f0 values, shown in Fig. 2, were extracted using RAPT (Talkin, 1995), a commonly used f0 detection algorithm, used in Qian et al. (2007)'s Cantonese suprataone tonal recognizer. Speaker-specific pitch floors and ceilings were set to the 1st and 99th quantiles minus or plus 30% of the range, respectively, a similar procedure to the pre-processing procedures in De Looze & Rauzy (2009); Evanini & Lai (2010).<sup>6</sup> Otherwise, the default parameter settings, including a 10ms frame shift were used. The first three and last frames were excluded because there were often large discontinuities between the estimated f0 for these frames and estimated f0 in the adjacent ones due to edge effects in the f0 detection algorithm, so there were a total of 67 f0 values, which were taken as the available f0 information in the intact condition. Unvoiced frames and frames with estimated f0 values resulting in large discontinuities were assigned f0 values using linear interpolation. To model the f0 information present in the degraded RESOLUTION conditions, the mean f0 was calculated over each unmasked region over frames falling within each of these regions. Thus, there were 6 f0 values estimated for each tritone in the 2-sample condition, one per unmasked region, and 9, 15, and 21 f0 values in the 3, 5, and 7-sample conditions, respectively. The f0 values were also log-transformed and then standardized as z-scores using speaker-specific means and standard deviations.

### 2.4. Participants

The participants were 39 native Cantonese speakers. There were 20 males (age  $18.9 \pm 1.8$  years) and 19 females (age  $21.9 \pm 1.9$  years). Participants were recruited from the local university student population at the City University of Hong Kong and at the University of California, Los Angeles and received cash compensation. All but three of the subjects (born/raised in Guangzhou and Shanwei, China) was born and/or raised in Hong Kong, China. Of the 10 participants tested in Los Angeles, all used Cantonese on a daily basis and had been in the United States for 3 to 8 years. No participants reported abnormal hearing.

### 2.5. Procedure

Participants were tested in sound-attenuated booths in the phonetics laboratories at the City University of Hong Kong and University of California, Los Angeles. The perception experiment was run in MATLAB using Psychophysics Toolbox extensions (Pelli, 1997; Brainard, 1997). Stimuli were played from an Echo Indigo IO sound card on a laptop over studio monitor headphones at a standardized, comfortable volume, and the responses and reaction times of the subjects measured from the

<sup>6</sup>The majority of the f0 values were in the mid range since each tritone stimulus consisted of two mid-level tones (33), yielding a center-heavy distribution of f0 values; thus, we could not use less extreme quantiles as in De Looze & Rauzy (2009); Evanini & Lai (2010) because they resulted in severe compression of the estimated range.

1 onset of the stimulus were recorded.<sup>7</sup> The interstimulus interval  
2 was 3s.

3 Participants were told that the stimuli were extracted from  
4 sentences *let<sup>25/35</sup> yiu<sup>33</sup> wai<sup>33</sup> wai mat<sup>3</sup> geng<sup>33</sup>* ‘You want  
5 NAME to clean the mirror,’ and they were given a sheet of pa-  
6 per with orthographic characters which showed what stimuli  
7 was being played, and what word they were to identify: *wai<sup>33</sup>*  
8 *mat<sup>3</sup>*. The stimuli were blocked by temporal resolution;  
9 block order was pseudorandomized to be roughly uniformly  
10 distributed over sampling resolution condition across partici-  
11 pants to average over learning effects between blocks, and stim-  
12 uli were randomized within blocks. Each block contained 90  
13 stimuli, 18 from each of the 5 speakers, with 3 distinct repeti-  
14 tions per speaker per target tone. The task of the participants  
15 was to lexically identify the target syllable in each stimulus by  
16 a keyboard press of one of six keys labeled with the characters  
17 for the minimal tone set over *wai*. Participants were asked to  
18 respond as quickly and accurately as possible and told that they  
19 would be timed.

## 21 22 2.6. Data analysis

23 Statistical analysis was performed in R (R Core Team,  
24 2014), and graphics were created using the ggplot2 package  
25 (Wickham, 2009). Tone identification accuracy was tested  
26 against at-chance levels (1/6) in each resolution condition by  
27 checking if 1/6 was contained in the 95% high posterior den-  
28 sity interval of Bayesian estimates of the mean accuracy across  
29 tones (Kruschke, 2013). For the human perceptual data analy-  
30 sis, all listeners and items were included in the analyses. While  
31 all listeners performed at above chance levels in the intact con-  
32 dition overall, not all listeners performed above chance levels  
33 for each individual tone, and in addition, there were three items  
34 that were not identified at above chance levels.

## 35 36 37 2.6.1. Machine classification

38 For insight into the human perception results, the acoustic  
39 separability of the different tones in the stimulus set was as-  
40 sessed using the acoustic feature extraction described in §2.3.  
41 These acoustic features were used in machine classification  
42 of the tone stimuli, using support vector machine classifiers  
43 (SVMs) (Vapnik, 1995; Cortes & Vapnik, 1995; Burges, 1998).  
44 SVMs are well-understood and widely used in machine learn-  
45 ing and have been popular for automatic tonal recognition,  
46 e.g. Peng et al. (2004); Levow (2005); Peng & Wang (2005);  
47 Wang & Levow (2008); Wang et al. (2009); Chen et al. (2014);  
48 Lam (2014). For an intuitive explanation of how SVMs work,  
49 see Appendix B.

50 Because the SVM algorithm involves calculating Euclidean  
51 distances in the parameter space, it is necessary to scale the  
52 data, so that parameters with a greater range do not dominate  
53 the direction of the optimal separating hyperplane relative to  
54 parameters with a smaller range. Thus, the f0 data was log

55 transformed and then z-score standardized, following Levow  
56 (2006, §2.3). (This preprocessing step was also similar to that  
57 of Peng & Wang (2005), which used a log-transformed 5-level  
58 normalization.) An SVM classifier was built for each of the  
59 2-, 3-, 5-, and 7-sample conditions and 10 ms frameshift “in-  
60 tact” condition, as well as three 1-sample conditions, using ei-  
61 ther the first, second, or third sample of each syllable from the  
62 3-sample condition, “1-sample-initial”, “1-sample-medial” and  
63 “1-sample-final”. These 1-sample conditions are included in  
64 the discussion only when explicitly mentioned.

65 Linear SVMs were implemented with LIBSVM  
(Chang & Lin, 2001) in R’s e1071 package (Meyer et al.,  
2012). The functions that linear SVMs use to classify data  
are linear combinations of the features, e.g., for a resolution  
condition with  $n$  f0 samples in total over the three syllables,  
 $\alpha_1 \cdot f0_1 + \alpha_2 \cdot f0_2 + \dots + \alpha_n \cdot f0_n$ , where, for all  $i$ ,  $\alpha_i$  is a  
constant. Thus, the absolute value of the feature weight can  
be used as a measure of a feature’s relative importance for the  
classifier (Guyon & Elisseeff, 2003).

The 6-way Cantonese tone classification problem was de-  
composed as  $\binom{6}{2} = 15$  binary classification sub-problems, e.g.  
T55 vs. T25, and the tone category receiving the most votes  
over all the sub-problems was selected as the classification de-  
cision. For each sampling resolution condition, the data was  
partitioned into 5 folds, one fold per speaker, for 5-fold cross-  
validation. Rotating across the folds, a single fold (18 tritones,  
1 speaker) was used as training data, and the remaining four  
folds ( $4 \times 18 = 72$  tritones, 4 speakers) were used as test data.  
All classification results, unless otherwise indicated, were aver-  
aged across the results from the 5 rotations, and standard error  
for classification accuracy was calculated from the variance of  
the accuracy over the 5 folds. Tonal classification accuracy was  
analyzed with logistic regression as described in §2.6.2.

## 2.6.2. Logistic models for tonal identification responses

The probability of correctness of tonal identification was an-  
alyzed using mixed effects logistic regression implemented by  
the lme4 package of Bates et al. (2014). Since exploratory data  
analysis showed that the effect of RESOLUTION on tonal iden-  
tification accuracy varied by tone, (see Fig. 3), we analyzed  
separate models of identification accuracy for each of the six  
tones. The logistic models included the fixed effect of RESOLU-  
TION, which was specified in two ways: (1) as contrasts between  
each degraded RESOLUTION condition and the intact condition,  
and (2) as forward difference contrasts between each successive  
RESOLUTION condition, e.g. between the 3-sample and 2-sample  
condition. Logistic models were also used to analyze the effect  
of RESOLUTION on the probability of response of a tonal cate-  
gory for the least accurately identified tones and the tones they  
were most confused with: responses of T22 and T23 for T33,  
responses of T25 and T22 for T23, responses of T23 for T22,  
and responses of T23 for T25.

To avoid anticonservativity, the random effects structure was  
chosen to be the maximal random effects structure justified by  
the experimental design that led to convergence (Barr et al.,  
2013); this procedure resulted in the inclusion of random in-  
tercepts by (stimulus) speaker and listener, as well as random

59 60 61 62 63 64 65 <sup>7</sup>Reaction times were measured but not further reported here since the 2-  
and 3-sample conditions had longer stimuli than the other conditions and since  
no significant effects were found for temporal resolution.

1 slopes for RESOLUTION by listener except when the model did not  
2 not converge with the random slope included. For SVM data,  
3 we included random intercepts by speaker and fold, as well as  
4 random slopes for RESOLUTION by speaker when the model con-  
5 verged with them. Listed p-values for fixed-effects coefficients  
6 are from Wald z-statistics. Significance was determined at an  
7 alpha level of 0.05. (Full regression coefficient values and p-  
8 values are given in Supplementary Materials, §5).

### 10 2.6.3. Analysis of confusion matrices

11 For visualizing confusion patterns in tonal responses, the  
12 confusion matrices from human perception and machine clas-  
13 sification were also analyzed using the well-studied similar-  
14 ity choice model (Luce, 1963; Nosofsky, 1990), following  
15 Silbert (2014). This model allowed us to partition the un-  
16 derlying source of the tonal responses into the distinct con-  
17 tributions of similarity between stimuli and bias in responses.  
18 The calculated similarity matrices were used as input to non-  
19 metric 2-D multidimensional scaling (MDS) implemented us-  
20 ing smacof (de Leeuw & Mair (2009)), and average linkage  
21 hierarchical clustering (James et al., 2013, p. 390-396) using  
22 hclust (R Core Team, 2014). For the MDS solutions, stress  
23 was  $5e^{-3}$  or below for all resolution conditions except for the  
24 intact condition, where stress was 0.01.

## 25 3. Results

26 This section describes results on the effects of temporal res-  
27 olution on acoustic classification and perceptual identification  
28 of Cantonese tones. The discussion of these effects is bro-  
29 ken into three sections addressing our hypotheses: (a) results  
30 that provide evidence for the sufficiency of coarse temporal res-  
31 olution in acoustic classification and perceptual identification  
32 ( $\S 3.1$ ), (b) results that show that coarser resolution affects dif-  
33 ferent tones differentially ( $\S 3.2$ ), and (c) results that compare  
34 the informativity of the syllable offset vs. the midpoint and the  
35 onset ( $\S 3.3$ ).

### 36 3.1. The sufficiency of coarse temporal resolution

37 The evidence for the sufficiency of coarse temporal res-  
38 olution comes from three results: (a) tonal identification was  
39 well above chance for humans and machine even under coarse  
40 sampling ( $\S 3.1.1$ ), (b) decreasing temporal resolution had quite  
41 limited effects on identification ( $\S 3.1.2$ ), and (c) for both hu-  
42 mans and machine, the multidimensional scaling and hierar-  
43 chical clustering solutions were very similar across resolutions  
44 ( $\S 3.1.3$ ).

#### 45 3.1.1. High identification accuracy under coarse sampling

46 Tonal identification was well above chance ( $1/6 = 16.67\%$ )  
47 for every resolution condition for humans and machines, even  
48 when humans had less than a quarter of the speech signal avail-  
49 able with 2 samples/syllable, and when SVM input consisted  
50 of only 1 sample/syllable. For humans, it ranged from 67.46%

51 (SE 2.91) in the intact condition to 60.51% (SE = 2.41) in the 3-  
52 sample condition and 52.54% (SE 2.90) in the 2-sample condi-  
53 tion. With only 3 samples per syllable, human identification accu-  
54 racy for T55 was 87%, 71% for T21, 70% for T25, and 62%  
55 for T33. SVM classification accuracy of z-scores ranged from  
56 65.00% (SE 2.86) in the intact condition (with 67 f0 values per  
57 exemplar) to 63.61% (SE 3.16) for the 2-sample condition. Accu-  
58 racy was still 62.78% (SE 2.96) even in the 1-sample-medial  
59 conditions: classification accuracy was 100% for T55, 72% for  
60 T25, and 67% for T21, and 58% for T33 and T22.

61 Figures 3 displays the relative frequency of responses for  
62 each individual tone in human perception, conditioned on RES-  
63 OLUTION. Each panel shows the distribution of responses for a  
64 given tone, indicated by the label at the top of the panel. The  
65 response tone is indicated by text labels at each data point, e.g.  
66 the identification accuracy for T55 can be read off from the se-  
67 ries of points marked as “55” in the panel labeled “55”. Figure 4  
68 is like Figure 3, but for SVM classification of z-score standard-  
69 ized log-transformed f0 values. These figures show that both  
70 humans and SVMs recognized T55 with the highest accuracy.  
71 However, while T55 accuracy for humans was around 85%, it  
72 was 100% across nearly all resolutions for the SVMs. Also  
73 like humans, the SVMs recognized T22 and T23 with around  
74 30-50% accuracy, the lowest out of all tones, and T25 and T33  
75 with 70-75% accuracy. SVM accuracy for T21 recognition was  
76 around 53%, though—much lower than human accuracy for  
77 T21, which was 70% and above.

78 The reader may wonder why accuracy in the intact condition  
79 was so low for both humans and SVMs. From Figures 3 and  
80 4, it is clear that for both humans and SVMs, the poor per-  
81 formance in the intact condition can be traced to errors on T22 and  
82 T23, and for SVMs, on T21, too. Identification accuracy was  
83 43% for T22 and T23 in the intact condition for humans, but  
84 77% for the other four tones. Moreover, identification accuracy  
85 for the other four tones dropped only to 73% on average in the  
86 3-sample condition and 65% on average in the 2-sample con-  
87 dition.<sup>8</sup> SVM classification accuracy in the intact condition for  
88 T55, T25, and T33 was 81% on average, but 49% for the other  
89 three tones. Accuracies in the 2-sample condition were within  
90 2% of the intact ones. We return to discussion of consistently  
91 low accuracies for particular tones in §4.1.

#### 92 3.1.2. Limited drops in accuracy with decreasing resolution

93 Figures 3 and 4 also show that decreasing sampling resolu-  
94 tion resulted in only quite circumscribed drops in accuracy rela-  
95 tive to accuracy with the full speech signal.<sup>9</sup> There were little  
96 effects on human tonal identification accuracy down to 5 sam-  
97 ples/syllable, and pervasive detrimental effects only at 2 sam-  
98 ples/syllable.

99 Below, we present results from logistic mixed effects models  
100 comparing the odds of correct tone identification in degraded

<sup>8</sup>It is important that not only T33 but also three other tones showed relatively high accuracy in the intact condition and lack of sensitivity to sampling resolution, since it's possible that the occurrence of two consecutive *wai*<sup>33</sup> syllables for the T33 stimuli may have induced repetition-related effects in perception.

<sup>9</sup>Confusion matrices with numerical values are given in the Supporting Ma-  
105 terials in §2.

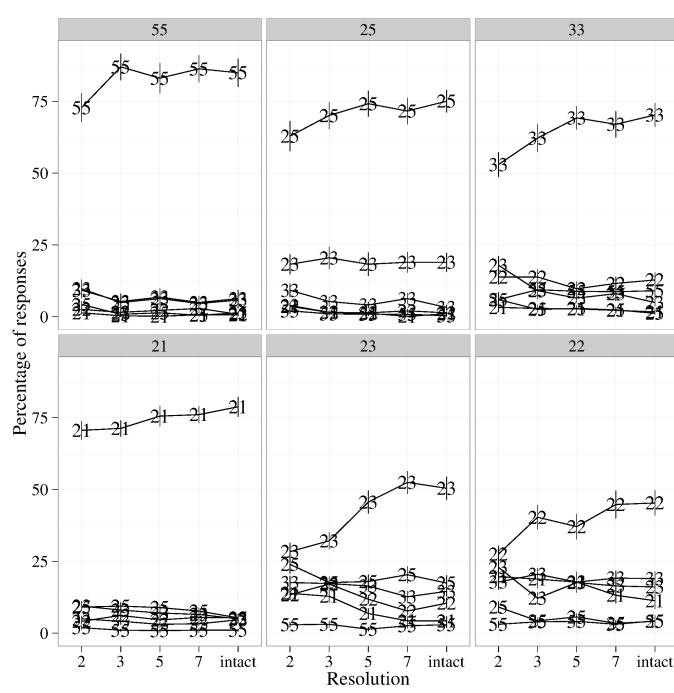


Figure 3: Cantonese native listeners' tonal identification response frequencies for each of the six tones, conditioned on RESOLUTION. Accuracy for all tones except Tone T22 was significantly lower than in the intact condition only in the 2- and 3-sample conditions. T23 and T22 were identified with strikingly lower accuracy overall than the other tones were. Error bars show  $\pm 1\text{SE}$  over participants.

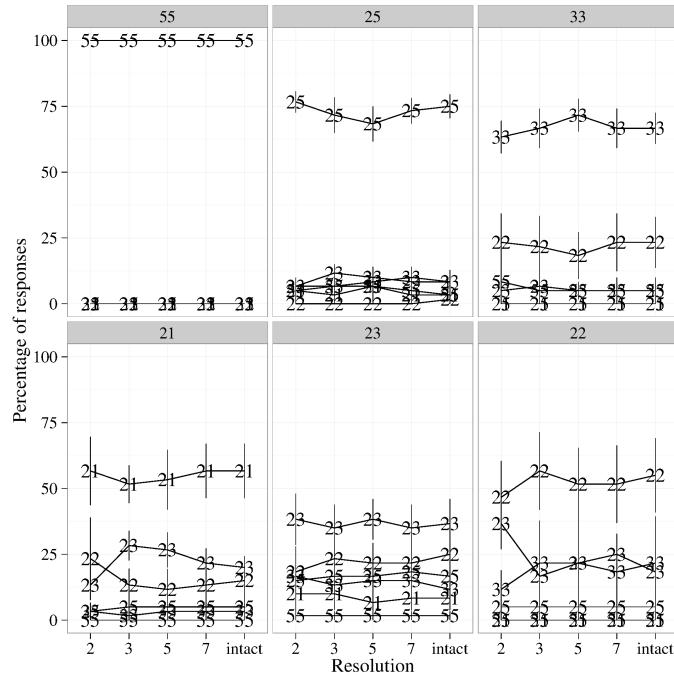


Figure 4: Support vector machine classification response frequencies for each of the six tones, conditioned on RESOLUTION. Like in human perception, T55 was recognized with highest accuracy and T22 and T23 with the lowest accuracy. Unlike in human perception, T21 was also identified with relatively low accuracy. Error bars show  $\pm 1\text{SE}$  over folds.

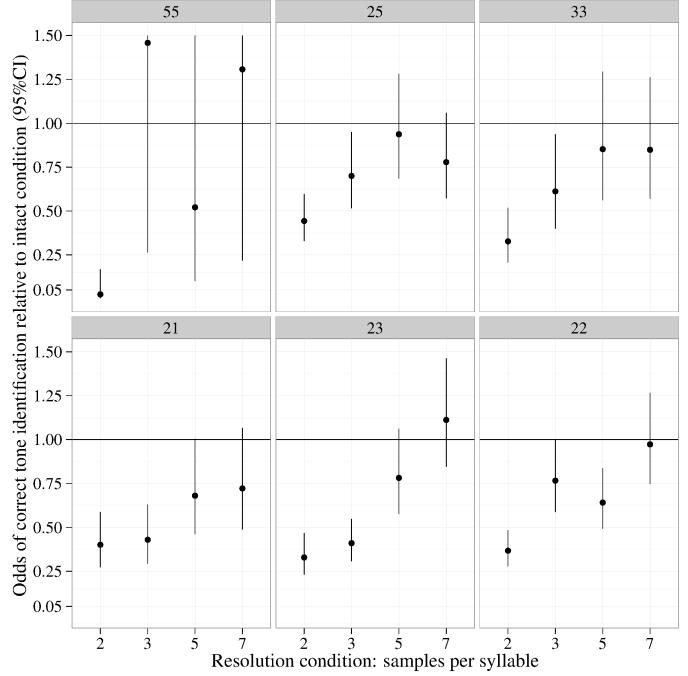


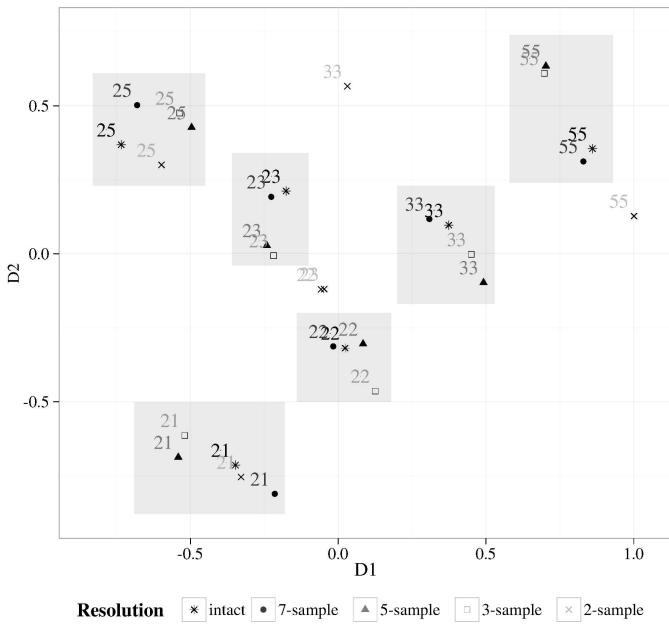
Figure 5: Ratio of odds of correct tone identification by humans for RESOLUTION conditions to odds of correct tone identification by humans for the “intact” condition, with estimated 95% confidence intervals. Confidence intervals estimated using standard errors of fixed effects.

RESOLUTION conditions to the odds in the intact condition. Results for humans are summarized in Figure 5. Each panel shows (for the tone labeling the top of the panel) the ratio of the odds of correct tone identification for the degraded RESOLUTION conditions to the odds of correct tone identification for the “intact” condition, with estimated 95% confidence intervals; numeric values are given in §5.1 in the Supplementary Materials. For cases where the confidence interval includes 1, there was no significant difference between the odds of correct tone identification and the odds in the “intact” condition.

For T55, only the 2-sample condition showed a significant difference in the probability of correct tone identification. For T33, T21, and T23, both the 2-sample and 3-sample conditions showed a significant difference. For T22, 2-, 3-, and 5-sample conditions showed a significant difference. Thus, for all tones except T22, there were significant drops in tonal identification accuracy in the degraded resolution conditions only for the 2- and 3-sample conditions. In fact, for T55, T25, T33, and T22, the 95% confidence interval for the odds of correct tone identification relative to in the intact condition was 0.94 or above, even in the 3-sample condition. Temporal resolution therefore had pervasive detrimental effects on human tonal identification accuracy only in the 2-sample condition.

The effect of decreasing resolution in SVM classification was even weaker than in human perception (See §5 in Supplementary Materials for SVM logistic regression results.). In logistic models for SVM classification, we did not include T55 since it was recognized with near perfect accuracy across resolutions. We found significant effects for RESOLUTION only for T25: for

1 z-scores, there was a significant decrease in accuracy for T25  
 2 in the 5-sample condition ( $\beta = -0.53$  (SE 0.23),  $z = -2.28$ ,  $p = 0.02$ ). Between successive resolution conditions, only T22  
 3 showed any significant difference in odds of correct tone identification, and this drop in odds occurred between the 3- and  
 4 2-sample conditions.



32 Figure 6: 2-D multidimensional scaling solutions for similarity matrices de-  
 33 rived from human tonal perception responses for each RESOLUTION condition.  
 34 Regions enclosing a cluster of identical tones are highlighted with grey boxes.  
 35 The 2-sample MDS space is quite different from all others.

### 3.1.3. Common acoustic and perceptual spaces across resolutions

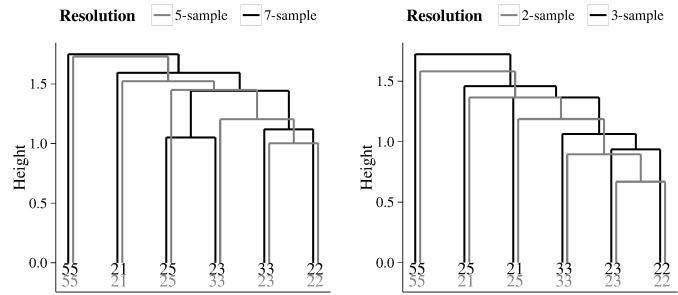
40 More evidence for a limited effect of RESOLUTION on tonal  
 41 identification comes from visualizations of the perceptual and  
 42 acoustic spaces underlying tonal identification. Visualizations  
 43 of the effects of RESOLUTION on the similarity matrices derived  
 44 from the human responses are summarized in an overlay plot  
 45 of 2-D MDS solutions for the different resolutions (Figure 6)  
 46 and in dendrograms from average linkage hierarchical clustering  
 47 (Figure 7). The gross distribution of tones in 2-D spaces  
 48 derived from non-metric multidimensional scaling of human  
 49 confusion data changed very little down to the 3-sample con-  
 50 dition. T55 was quite isolated at the periphery of the space for  
 51 all resolutions. The tones T25 and T21 also consistently ap-  
 52 peared rather aloof at the periphery of the MDS space across  
 53 resolutions, while T33, T23, and T22 crowded together in the  
 54 center of the space.

55 This same periphery-center distinction was mirrored in the  
 56 dendrograms in Figure 7. The height at which tone categories  
 57 fused in the dendrogram indicates how similar they were: tone  
 58 categories that fused near the top were very dissimilar from one  
 59 another, while tone categories that fused near the bottom were

60 very similar to one another. Thus, since T55 fused at the top  
 61 of every dendrogram in Figure 7, it was the tone that was the  
 62 most perceptually dissimilar from all the other tones in every  
 63 RESOLUTION condition. Tones T25 and T21 also fused near the  
 64 top of every dendrogram, while the highly mutually confusable  
 65 T33, T23, and T22 tones fused at the bottom of every dendro-  
 66 gram. The same was true for every dendrogram for SVM clas-  
 67 sification. (See §3 in Supplementary Materials for SVM z-score  
 68 dendrograms.)

69 The only major shift in the perceptual space down to the 3-  
 70 sample condition was that T25 rose out of a bottom cluster with  
 71 T23 from the 7- to the 5-sample dendograms (Figure 7, left). This  
 72 was because the accuracy of correct T23 identification be-  
 73 tween those conditions dropped by 6.8%, a significant decrease  
 74 in odds of correct identification ( $\beta = -0.35$  (SE 0.15),  $z = -2.30$ ,  
 75  $p = 0.021$ ). Between the 7- and 5-sample conditions, T23, T22,  
 76 and T33 retracted downwards into the interior of the perceptual  
 77 space, away from T25, as shown in Figure 6.

78 The distribution and location of tones in 2-D multidimen-  
 79 sional scaling solutions for SVM confusion data was almost  
 80 insensitive to RESOLUTION down to 1 sample, especially if only  
 81 the 1-sample-medial condition was included, see Supple-  
 82 mentary Material §4. Points for any particular tone were tightly  
 83 clustered together across resolutions. Like the human MDS  
 84 space, the SVM space had T33, T23, and T22 right next to one  
 85 another, and T55 isolated on the periphery. The one major dif-  
 86 ference from the human space was that the SVM space has T21  
 87 and T25 on the same side of the periphery next to one another.  
 88 This indicated that T21 and T25 were much more similar to one  
 89 another for the SVMs than for humans.



90 Figure 7: Dendrograms from average linkage hierarchical clustering of simi-  
 91 larity matrices. Similarity matrices were calculated from the similarity choice  
 92 model applied to human perception data.

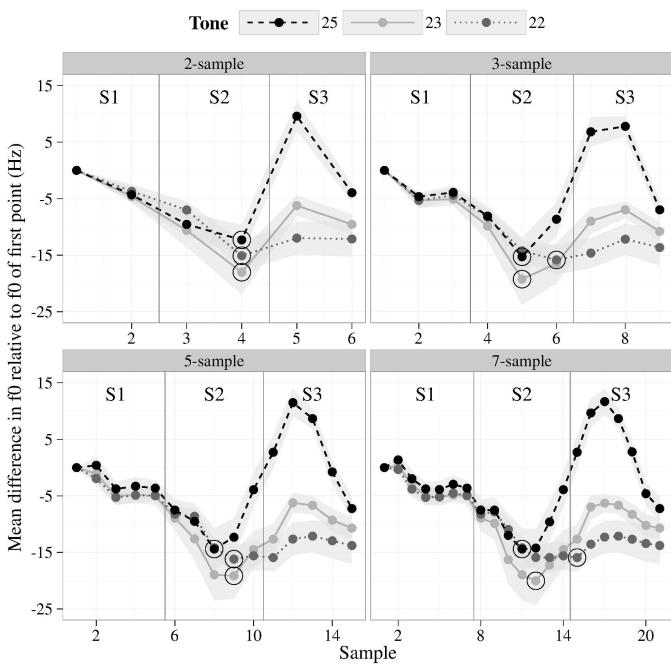
### 3.2. The insufficiency of coarse sampling resolution for partic- 93 ular tones

94 Close examination of confusion patterns, SVM feature  
 95 weights, and f0 contours suggested that loss of detailed f0 tem-  
 96 poral alignment and contour shape information with decreased  
 97 resolution was a major cause of decreasing accuracy in tonal  
 98 identification for identifying a subset of the tones. We describe  
 99 three examples of evidence for this below: (a) the collapse of  
 100 different f0 minima alignments between rising contours under  
 101 coarse resolution (§3.2.1), (b) the role of loss of information

1 about contour shape in the drop in T22 identification accuracy  
 2 between the 3- and 2-sample conditions (§3.2.2), and (c) the  
 3 consequences of undersampling of the f0 extremum for T23 vs.  
 4 T25 and T21 vs. T23 classification (§3.2.3).

### 5 3.2.1. Collapse of differing alignment of f0 minima in rising 6 contours

7 The circled minima of the f0 contours of T25, T23, and T22  
 8 in Figure 8 highlight the effect of decreasing resolution on the  
 9 discriminability of rising f0 contours. As resolution decreased,  
 10 the differences in the timing of the turning point in the f0 con-  
 11 tours vanished. The f0 minimum of T25 was shifted later, to  
 12 coincide with the minima of T23 and T22; the broad f0 valley  
 13 of T22 was shifted earlier and sharpened, reducing T22’s dis-  
 14 distinctness from T23; the depth of the f0 minimum for T23 was  
 15 raised, reducing the distinctness between the T23 minimum and  
 16 that of the other two tones. In the 2-sample condition, the cir-  
 17 cled points for all three tones appear in a vertical line at sample  
 18 number 4: there’s no more distinction between the alignment  
 19 of the f0 minima of the tones. The collapse of f0 minima align-  
 20 ment and shape resulted in increased confusability of T22 and  
 21 T23 and of T23 and T25—this is discussed in §3.2.2, which  
 22 follows.



49 Figure 8: Effects of sampling resolution on resolving turning points for T25,  
 50 T23, and T22 stimuli. Each panel displays the f0 information present for the  
 51 RESOLUTION indicated at the top of the panel, with linear interpolation between  
 52 points. For each individual stimulus, the f0 value of the first point was sub-  
 53 tracted from each f0 point. Each tone’s f0 contour is averaged over all stimuli  
 54 for the tone, and ribbons show  $\pm 1\text{SE}$ . A circled point indicates the global min-  
 55 imum in the aggregate contour.

### 56 3.2.2. Contour shape and confusability of T22 and T23

57 Sampling f0 from only syllable onsets and offsets (the 2-  
 58 sample condition) significantly reduced tonal identification ac-

curacy from sampling medially as well as at onset and offset (the 3-sample condition). In the 3-sample condition, only T21 and T23 saw identification accuracy drop relative to conditions with finer sampling for humans (see Figure 5). But sampling resolution had pervasive detrimental effects for every tone in human tonal identification only in the 2-sample condition.

Confusion patterns also shifted greatly from the 3- to 2-sample condition. This shift is most striking in the 2-sample MDS space, where the periphery-center divide in tonal distribution was destroyed. T22, T23, and T33 all shifted away from T21 and closer to T25. However, while T22 and T23 remained in the center, nearly coinciding in space, T33 was shifted out into the periphery, nearest to T25. The high confusability between T22 and T23 was also reflected in the 2-sample dendrogram (Figure 7, right), where the height of fusion for T22 and T23 was by far the lowest. Moreover, the odds of correct identification of T22 significantly decreased from the 3- to 2-sample condition for humans from 40.34% to 27.52% ( $\beta = -0.58$  (SE 0.21),  $z = -2.84$ ,  $p = 4.6e^{-3}$ ) and for SVMs from 56.67% to 46.67% ( $\beta = -2.9$  (SE 1.06),  $z = -2.72$ ,  $p = 6.5e^{-3}$ ). This was the only significant change for SVM in odds of correct tone identification between successive resolution conditions. At the same time, the odds of a T23 response to T22 significantly increased for humans by 10% ( $\beta = -0.58$  (SE 0.21),  $z = -2.84$ ,  $p = 4.6e^{-3}$ ) and for SVMs by 20% ( $\beta = 2.60$  (SE 1.05),  $z = 2.48$ ,  $p = 0.013$ ).

Why were T22 and T23 increasingly confusable as resolution dropped? Figure 9 suggests that the source of the confusion was the collapse of distinctions in the alignment and shapes of the T23 and T22 valleys. The plots of SVM feature weights in Figure 9 show that there were two time intervals of particular importance in acoustic discrimination of T22 and T23: the descent to the f0 valleys in syllable 2 and the f0 peaks in syllable 3. In the 2-sample condition (top center panel), the relative importance of f0 information from the descent to the valleys (samples 3 and 4) plummeted.

A comparison of f0 contours across resolutions shows that f0 information from the 2-sample condition missed the difference in slope of descent between the sharp fall of T23 and the flat basin of 22. The lack of this information may have been the source of the increased confusability of T22 for T23 between 3- and 2-sample conditions in humans and SVMs (see §3.1.2).

### 3.2.3. Undersampling the T25 and T23 f0 peaks

For the classification between T25 and T23 (see Figure 10), it was the time interval where the peaks occurred that had the most important f0 information.<sup>10</sup> As the signal was degraded, the number of samples where T25 and T23 were at their f0 peaks dropped from many in the “intact” condition to 2 samples

<sup>10</sup>Like in Figure 9, there is also a secondary peak in the SVM feature weight timecourse in Figure 10. The secondary peak precedes the rise in the middle of syllable 2 and is missing in the 2-sample condition. However, this is actually because of a missing valley in the SVM feature weight timecourse between the two peaks: the 2-sample condition misses any f0 samples between the valleys and peaks of the T23 and T25 contours. In that time interval, the T23 rise can be quite steep, “catching up” to the level of the T25 rise after having dropped to a lower valley than T25.

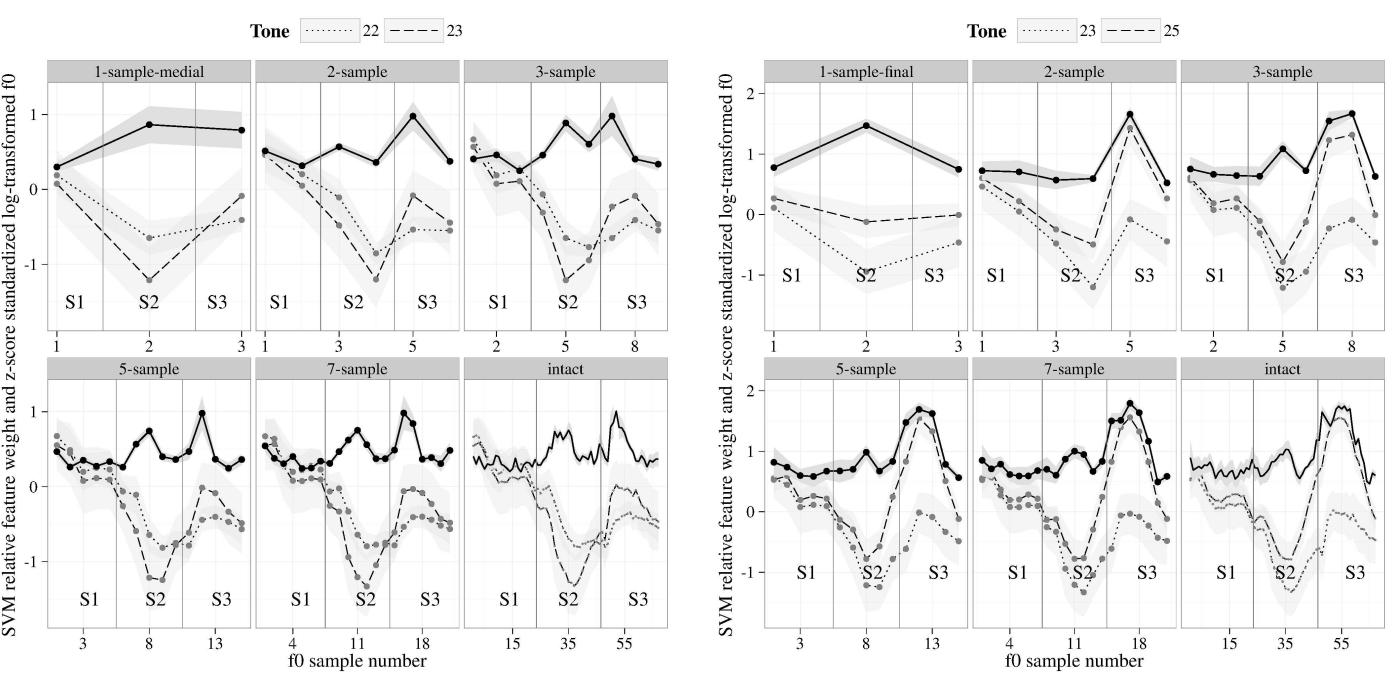


Figure 9: Averaged z-score standardized log-transformed f0 contours (in grey) for T23 and T22 with time-aligned SVM feature weight magnitudes (in black) for different resolution conditions. Ribbons show  $\pm 1\text{SE}$ .

in the 3-sample condition and just 1 sample in the 2-sample condition. Correspondingly, the peak in SVM feature weight coincident with the f0 peaks sharpened from a plateau to a point as the signal is degraded. Having just a single sample to detect the differences in peak height between T25 and T23 in the 2-sample condition may have been behind the significant 5% increase in T25 responses for T23 stimuli from the 3- to 2-sample condition ( $\beta = 0.44$  (SE 0.15),  $z = -2.86$ ,  $p = 4.3\text{e}^{-3}$ ). There were no significant differences in odds of a T25 response for T23 between any other successive resolution conditions. For SVMs, there were no significant differences in odds of a T25 response for T23 at all.

### 3.3. Comparison of syllable onset, midpoint, and offset as locations for sampling

Comparing tonal classification by machine under sampling just once per syllable at syllable onset, midpoint, and offset, we found that: sampling at just the onset resulted in missing the valleys of T25 and T21 (§3.3.1); sampling at just the midpoint made rises and falls very confusable (§3.3.2), and sampling at just offset was severely detrimental to T25 identification (§3.3.3).

#### 3.3.1. Sampling at syllable onsets misses valleys

Sampling at just the syllable onset also resulted in poor tonal classification. While T55 identification accuracy was perfect for the other 1-sample conditions, it was only 70.00% (SE 5.65) when sampling only at syllable onset. This is because T55 became confusable with T25, since sampling at the onset results in missing the dip before the rise of T25. Moreover, since sampling

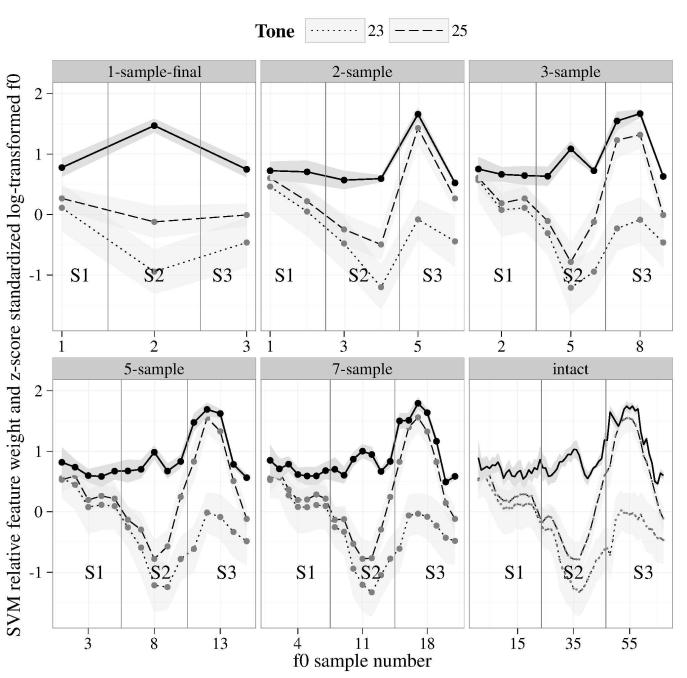


Figure 10: Averaged z-score standardized log-transformed f0 contours (in grey) for T23 and T25 with time-aligned, scaled SVM feature weight magnitudes (in black) for different resolution conditions. The points were linearly interpolated. SVM weight magnitudes within a panel were scaled by a constant to make their relative differences easier to discern. Vertical lines within each panel delineate syllable boundaries. Ribbons show  $\pm 1\text{SE}$ .

at the onset results in missing the valley of T21, T21 was particularly confusable with T23 and T22 and classified correctly 36.67% (SE 11.96).

#### 3.3.2. Rises and falls: low f0 at syllable midpoint

Sampling at just syllable midpoint resulted in shrinking the distinction between rises and falls. Rises and falls were not distinguished by L vs. H targets at the syllable midpoint. The f0 contours in Figure 11 show that the T21 fall did not have a high target near syllable midpoint—it was high only at syllable onset. Instead, the T21 fall had an f0 minimum between the syllable midpoint and offset. Moreover, with decreasing sampling resolution, this f0 minimum appeared to shift to the right.

The location of the f0 minimum for T21 was also just where the f0 contours for T25, T23, and T22 achieved their f0 minima as well. This can be seen from Figure 8, which shows f0 contours for T25, T23, and T22 aggregated over all stimuli in the experiment by tonal category, for all degraded resolution conditions. (For easier interpretability, the f0 contours were standardized by subtracting off the first f0 value in the contour (the “anchor”) from all other f0 values rather than shown as z-scores.) The f0 contours for these three tones descend to their f0 minima between the syllable 2 midpoint and offset, before beginning their ascents.

T21 was in fact most confusable with rising tones in human and SVM classification. In the human perception 2-sample condition, the frequency of responses of T23 and T25 to T21 was 9% for each rise, at least twice as much as responses for any

other tone. Sampling only at the syllable midpoint, SVM identification of T23 was only 22%, with a 22% rate of T21 responses. In contrast, sampling at either the syllable onset or offset resulted in only 3-5% T21 responses. From the top left panel in Figure 11, we can see that sampling at only syllable midpoint shrunk the difference in f0 heights at the f0 minima of T21 and T23.

### 3.3.3. Insufficiency of syllable offsets for identifying T25

Sampling at just syllable offset resulted in shrinking the distinction between T25 and level tones. Sampling f0 from only syllable offsets was severely detrimental to identification of T25 in SVM classification. With just 1 sample per syllable, SVM classification accuracy for T25 was just 25.00% (SE 6.45) at syllable offset, compared to 53.33% (SE 10.07) at syllable onset and 71.67% (SE 4.25) at syllable midpoint. In the 1-sample-offset condition, T25 was classified as the level tones T22 and T23 55% of the time. In comparison, T25 was mistaken for T22 and T23 only 15% of the time for the other 1-sample conditions.

The reason for the insufficiency of f0 information from the syllable offset is clear from Figure 10. This is a timecourse plot of z-score standardized log-transformed f0 points overlaid with scaled absolute values of SVM feature weights for each f0 sample for T25 vs. T23 classification. The f0 contours are aggregated over all stimuli in the experiment by tonal category. Figure 10 shows that the T23 and T25 rises both exhibited peak delay, with their peaks occurring in the syllable following the one that they were associated to (in syllable 3). Moreover, the f0 values most important in machine classification of these tones occurred at these peaks in the first half of the following syllable—not at syllable offset. When f0 information was sampled from just the offset (1-sample-final condition, top left panel), the ascent and descent of the T25 peak was missed almost entirely. This is why T25 was highly confusable with level tones in the 1-sample-final condition. The fidelity of the shape of the T23 contour with just syllable offsets was higher since f0 coming out of the f0 minimum was still relatively low at syllable offset in the 2nd syllable.

## 4. Discussion

This study provided evidence for both the sufficiency and insufficiency of coarse sampling resolution in the acoustic classification and perceptual identification of tones. The evidence came from native speaker perceptual identification of Cantonese tones degraded by interrupting noise, as well as acoustic classification of the tones by support vector machines from fewer and fewer f0 data points in the input.

As we hypothesized, overall, tonal identification accuracy by humans and machine was well above chance for all temporal resolution conditions:

1. Tonal identification accuracy above 70% was maintained for T55 and T25, even with just 3 samples/syllable for humans and 1 for SVMs, and accuracy for each tone in every resolution condition was well above chance.

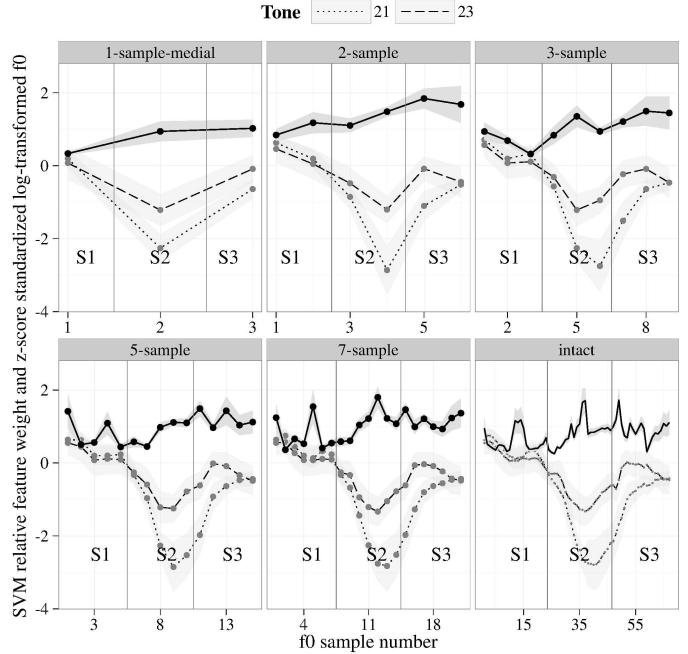


Figure 11: Averaged z-score standardized log-transformed f0 contours (in grey) for T21 and T23 with time-aligned SVM feature weight magnitudes (in black) for different resolution conditions. Both tones had f0 minima between the syllable midpoint and offset.

2. Decreasing sampling resolution resulted in only limited decreases in identification accuracy—pervasive deterioration of identification accuracy occurred only at 2 samples/syllable for humans, and SVM classification accuracy with 1 sample/syllable from the syllable midpoint was indistinguishable from accuracy with 67 f0 values/syllable.
3. Acoustic and perceptual spaces for tonal identification computed from multidimensional scaling and hierarchical clustering remained very similar across resolution conditions, except for the 2-sample condition in human perception.

We also hypothesized that the effect of decreasing temporal resolution would differentially affect different tones, and that identification of the two rising tones T25 and T23 would be especially affected. This was indeed the case:

1. Under coarse resolution, there was a collapse of differing alignment of f0 minima in the rising contours, making them increasingly confusable.
2. Dropping temporal resolution from 3 to 2 samples per syllable in humans, T22 and T23 became highly confusable; SVM results suggest that was because of loss in detail of contour shape at the f0 valleys.
3. SVM results suggested that increased confusability between T23 and T25 in humans from the 3- to 2-sample condition was due to loss of information about the shape of the f0 peaks.

Finally, we hypothesized that under sampling one f0 value per syllable, sampling at the syllable offset would facilitate

1 tonal classification more than sampling at the syllable onset or  
2 midpoint. However, we found instead that sampling at any one  
3 of these points resulted in poor identification of some subset of  
4 tones. In particular, sampling at the syllable offset resulted in  
5 very poor accuracy in identification of T25 by SVMs because  
6 of peak delay which placed the peak on the following syllable.  
7

8 In the remainder of this section, we further discuss results for  
9 each of the hypotheses: the sufficiency of coarse resolution in  
10 §4.1, as well as the insufficiency of coarse resolution for partic-  
11 ular tones and the effect of varying the location of sampling in  
12 §4.2.

#### 13 4.1. The sufficiency of coarse resolution

14 Our positive results for the sufficiency of coarse resolution  
15 support the hypothesis of Chao (1930) and others—namely, that  
16 no more than a few f0 samples per syllable are needed to pro-  
17 vide distinctive features for tonal categories. The programming  
18 error in the stimuli noted earlier that introduced longer noise  
19 intervals at stimulus offset in the 2- and 3-sample conditions  
20 (27 and 7 ms, respectively) does not weaken the result that de-  
21 creasing resolution has a limited impact on tonal identification.  
22 The effect of the error could have only been in the direction  
23 of decreasing accuracy at the two lowest SAMPLING RESOLUTION  
24 conditions due to interference and/or memory effects.  
25

26 In addition, the significantly lower accuracy in the 2-sample  
27 condition compared to that in the intact condition may have  
28 been due in part to a lack of perceptual continuity caused by the  
29 long duration of the interrupting noise intervals for our partic-  
30 ular experiment design. (Recall that manipulation of the sam-  
31 pling resolution involved an increase of the duration of the noise  
32 interval as sampling resolution decreased.) In support of this  
33 conjecture, Dannenbring (1976) showed that in nonspeech, for  
34 pure tones of 250 ms in duration interrupted by white noise,  
35 the mean continuity threshold between perceived continuity and  
36 discontinuity due to the interrupting noise was around 80 and  
37 100ms for steady state tones and tone glides, respectively. This  
38 indicates that the 90-ms noise interval duration in the 2-sample  
39 condition may have been close to the auditory threshold for per-  
40 ceiving continuity in our stimuli, which were 241 ms in dura-  
41 tion.  
42

43 A more potentially serious challenge to the positive results  
44 for the sufficiency of coarse resolution is that the relative insen-  
45 sitivity of tonal identification accuracy to sampling resolution  
46 may have been due to a “floor effect”, i.e., already poor ac-  
47 curacy in the intact condition, which then left little room for  
48 further deterioration with decreasing sampling resolution. As  
49 discussed in §3.1.1, poor overall performance by humans and  
50 SVMs in the intact condition was due to errors on T22 and T23,  
51 and for SVMs, errors on T21 as well.  
52

53 We conjecture that the low accuracy for T22 and T23 in hu-  
54 man perception was due in part to tonal mergers in some of the  
55 listeners<sup>11</sup>(Bauer et al., 2003; Mok & Wong, 2010a,b). How-  
56 ever, there is also evidence that these tones were particularly  
57

58 confusable with other tones in the low pitch range when coar-  
59 ticulated with neighboring mid-level tones in the Tone T33 -  
60 Tone X - Tone 3 experimental stimuli. The strongest piece of  
61 evidence is that the SVMs performed just as poorly as the hu-  
62 mans for T22 and T23. The SVM classification assumed that all  
63 six tonal classes were equiprobable and received equal numbers  
64 of exemplars for each tone class from each speaker: there is no  
65 sense in which there were tonal mergers for the SVMs. This im-  
plies that poor accuracy for T22 and T23 in SVM classification  
must have been a property of their acoustics.

66 If T23 and T22 were confusable when flanked by mid-level  
67 tones, then what about Tone 21? We hypothesize that Tone  
68 T21 was identified with high accuracy unlike its close neighbors  
69 T22 and T23 not only because it occupied a lower part of the  
70 pitch range than them, but also because of creaky voice quality  
71 cues, since Yu & Lam (2011, 2014) showed that the presence  
72 of creaky voice cues can boost T21 identification accuracy in  
73 Cantonese tone perception. In support of this hypothesis, iden-  
74 tification accuracy for T21 was relatively high for humans, but  
75 low in machine classification. Also, the key difference between  
76 the MDS spaces for humans and SVMs was that T21 and T25  
77 were much closer in the SVM acoustic space.

78 Most of the unvoiced frames in the RAPT f0 extraction came  
79 in T21 stimuli, since Tone T21 realization frequently had non-  
80 modal phonation, low amplitude, and even intervals of silence.  
81 Acoustic feature extraction did not capture this, and while there  
82 are more sophisticated rules for estimating the pitch percept in  
83 the presence of nonmodal phonation for the human listeners  
84 than the simple linear interpolation used over voiceless frames  
85 that we used, cf. the aberrant pitch contours for T21 in Figure 2,  
86 the poor accuracy for T21 identification by machine nevertheless  
87 suggests that a parameterization of the speech signal that  
88 references voice quality, beyond f0, is needed for both higher  
89 classification accuracy of T21 by machine, and for modeling  
90 what humans are doing.

91 Finally, a potential concern for the relevance of the SVM re-  
92 sults for understanding human perception of tone is that pro-  
93 viding z-scores as input gave the machines an unreasonable  
94 amount of information since the calculation of z-scores assumes  
95 the ability to parse speaker identity from the signal, as well as  
96 some knowledge about each speaker’s pitch range. However,  
97 it could be the case that z-scores might be a rough proxy for  
98 some constellation of information for calculating pitch range  
99 available to listeners, such as absolute f0 and aspects of voice  
100 quality, as well as previous experience with many speakers  
101 (Bishop & Keating, 2012). Moreover, classification accuracy  
102 for anchored f0 values using information only internal to each  
103 stimulus for pitch range information (see Figure 8), was just as  
104 high, ranging from 65.28% (SD 10.44) in the intact condition  
105 to 63.61% (SD 9.99).

#### 106 4.2. The insufficiency of coarse resolution and the effect of 107 where samples are taken

108 Our negative results for the sufficiency of coarse resolution  
109 bear on previous work characterizing tone using 1 or 2 f0 sam-  
110 ples from particular locations in the syllable. First, the poor  
111 SVM classification results for T25 with 1 f0 sample at syllable

112 <sup>11</sup>For a discussion of potential tonal mergers in listeners and other possible  
113 reasons for low accuracy for T22 and T23, see Supplementary Materials §6.

1 offset provide a counterpoint to Khouw & Ciocca (2007)'s results that f0 change over the 6th and 7th out of 8 subsyllabic  
2 segments accounted for about 70% of the variance in a Cantonese tonal identification perception experiment of isolated  
3 monosyllables. While we did not use f0 change in the acoustic  
4 feature set, the bottom center panel in Figure 10 suggests that  
5 f0 change at syllable offset might not have fared any better than  
6 f0 for discriminating T25 and T23, since both rises have similar  
7 slopes at syllable offset. Our results on the poor separability of  
8 T25 at syllable offset are also interesting in the context of a body  
9 of previous work on tonal coarticulation in Cantonese (Li et al.,  
10 2002, 2004; Wong, 2006a) and other (South)east asian languages (Mandarin: (Xu, 1997), Thai: (Gandour et al., 1992),  
11 Vietnamese (Han & Kim, 1974)). This literature has reported  
12 that rightward (carryover) coarticulation is stronger than leftward (anticipatory) coarticulation, so that tones in connected  
13 speech might be maximally separated near the offset of the syllable.  
14 It seems that if there is significant peak delay, as is the case is here, then the syllable onset might also be a region of  
15 maximal separation for some tones. Moreover, there is likely  
16 to be peak delay under a wide span of speech rates, as rising  
17 contours always begin rising close to the syllable offset, even  
18 under slower speech rates (Xu, 1998, 2001; Wong, 2006b).  
19

20 Second, while rises and falls in connected speech might be  
21 distinguished by L vs. H targets at syllable midpoint in Thai  
22 (Zsiga & Nitisoroj, 2007), this does not appear to be the case in  
23 Cantonese. At least in the T3 - X - T3 context of our experimental stimuli, both rises and falls had f0 minima near syllable  
24 midpoint (Figures 8, 11)—in fact, the T21 fall's f0 minimum was  
25 rather close to syllable offset. While both T25 and T23 rises and  
26 the T21 fall could be described as having L targets between syllable  
27 midpoint and offset, the f0 minima for the three tones are,  
28 on average, at rather different heights. Thus, although rises and  
29 falls are not distinguished by different phonological targets, i.e.,  
30 H vs. L, they are distinguished by different phonetic targets.  
31

32 Third, although Barry & Blamey (2004) proposed a simple  
33 2-D space for Cantonese tone with the dimensions of f0 at syllable  
34 onset and offset, the high increase in confusability between  
35 T22 and T23 from 3- to 2-samples/syllable for both humans  
36 and SVMs suggests that something else not captured by f0 at  
37 syllable onset and offset also plays a role in Cantonese tonal  
38 perception. The SVM results (Figure 9) show that the relevant  
39 information lost may be the shape of the valleys in the f0 contours.  
40 The use of creak in Cantonese T21 perception also implies  
41 that the perceptual space for (native) Cantonese tone perception  
42 goes beyond a 2-D space of f0 samples, see Yu & Lam  
43 (2011, 2014).

44 Finally, our results suggesting a role of contour shape and  
45 temporal alignment in tone identification fit with the growing  
46 body of literature on contrastive tonal alignment, intonational  
47 perception, and auditory tonal processing in the brain that  
48 clearly show that shape and alignment in the f0 contour play an  
49 important role in tonal production and perception.

## 5. Conclusion

This study adds to the small amount of work that tests the effect of temporal resolution in the speech signal on tone perception. In line with previous work on Mandarin (Gottfried & Suiter, 1997; Lee et al., 2008, 2009; Lee, 2009), it shows that just a few samples per syllable are enough for both humans and support vector machines to classify Cantonese tones with reasonable accuracy, without much difference in performance from having the full speech signal available.

It is important to note two things, though. First, these results are about the sufficiency of coarse resolution for *off-line* tonal classification, and not on-line tone processing. Eye movements from on-line processing of Mandarin tones have shown that tonal perception is an incremental process (Shen et al., 2013), and as mentioned in §1.1, electroencephalographic studies of tone processing also clearly demonstrates that details of contour shape rather than just linear segments interpolating a few points are linguistically encoded in the definition of tonal categories (Chandrasekaran et al., 2007; Krishnan et al., 2009).

Secondly, a closer look at the results from our study shows that even for off-line classification, the identification of some tones fails under coarse sampling, and that *where* in the syllable the samples are taken greatly impacts how informative they are for tonal identification. This points to further exploration of adaptive sampling, where temporal resolution is not fixed, but may depend on the shape of the f0 contour or other properties of the speech signal; e.g. landmark approaches which extract different parameters from the speech signal on different time scales, determined by when inflection points and local extrema occur in the time course (Jansen & Niyogi, 2009), or approaches where sampling is dependent on the spectral stability of the speech signal (House, 1990, 2004a,b).

The importance of where samples are taken also shows that informative properties of the speech signal to identify a tone may extend to a time window beyond the syllable that the tone is associated to. This has long been established from studies on effects of preceding context on perception, such as Wong & Diehl (2003), which shows that the identification of a word uttered at a given f0 is completely determined by the f0 height of the preceding syllable. But here, we show that information we might consider “intrinsic” to the syllable actually drifts into the following syllable due to peak delay, so the following syllable is part of the domain of tonal realization rather than a source of external contextual information. Thus, in phonetic descriptions of tone, studies of tonal dispersion, and modeling the learning of tones in connected speech, in addition to careful consideration of what temporal resolution to use, we should also consider a time window for parameterization that is larger than just the syllable a tone is associated to.

On a final note, we emphasize that the role of temporal resolution in phonetic spaces is an issue that transcends which phonetic property or concept is under study. All the issues on temporal resolution of f0 contours we have raised here also apply to any other time-course properties in tonal spaces, e.g. spectral balance and amplitude. As a more distant case in point, automatic vowel formant extraction soft-

1 ware commonly extracts formants from just a single point  
 2 in the vowel (Evanini, 2011; Rosenfelder et al., 2011), de-  
 3 spite substantial evidence that formant trajectories are im-  
 4 portant for vowel identification—especially for English, with  
 5 its preponderance of diphthongal vowels (Strange et al., 1983;  
 6 Nearey & Assmann, 1986). Moreover, simulations for learn-  
 7 ing vowel categories from phonetic data have used input data  
 8 consisting of just formant values measured at a single point  
 9 (de Boer & Kuhl, 2003; Vallabha et al., 2007; Feldman et al.,  
 10 2013), as have studies of the typology of vowel dispersion  
 11 (Liljencrants & Lindblom, 1972; Becker-Kristal, 2010). There  
 12 is thus a vast body of literature on vowels that rests on the ten-  
 13 tuous assumption that the relevant temporal resolution for for-  
 14 mant values is a single sample per vowel.  
 15

16 In conclusion, we hope to have convinced the reader that tem-  
 17 poral resolution in the speech signal is not a finicky technical  
 18 detail or an engineering problem, but a fundamental issue for  
 19 phonetics that merits attention from phoneticians.  
 20

## Acknowledgements

23 Suppressed for reviewing process to preserve anonymity.  
 24

## Appendix A. Raw and transformed f0 ranges for speakers of experimental stimuli

29 The calculated raw and transformed f0 range of the stimuli  
 30 for each of the five speakers is given in Table A.1.  
 31

Speaker	f0 (Hz)	log f0	z-score
f4	[165.89,241.00]	[5.11,5.48]	[-3.35,2.35]
f3	[106.42,179.47]	[4.67,5.19]	[-5.78,1.83]
m6	[125.88,176.36]	[4.84,5.15]	[-2.92,3.21]
m1	[83.87,145.92]	[4.43,4.97]	[-3.48,1.84]
m5	[61.44,140.20]	[4.12,4.79]	[-5.08,3.60]

39 Table A.1: Speaker-specific f0 range in speech materials, measured in Hz, after  
 40 log-transformation, and after standardization of log f0 with respect to speaker  
 41 means and standard deviations. The speakers are ordered from highest to lowest  
 42 maximum f0, following the same order from top to bottom in the plot of f0  
 43 contours by speaker in Figure 2 in the paper.  
 44

## Appendix B. Background on support vector machines

49 We sketch a geometrical characterization of how support vec-  
 50 tor machines work for the binary case, e.g. for two tone classes,  
 51 following Bennett & Bredensteiner (2000). Call the two classes  
 52 Tone A and Tone B. Each stimulus is parameterized as a real-  
 53 valued  $p$ -dimensional vector and labeled as belonging to either  
 54 Tone A or B. Thus, the Tone A and B stimuli sets each comprise  
 55 a set of points in  $\mathbb{R}^p$ . The SVM algorithm is a way to determine  
 56 an optimal decision rule to assign a tone class label to a stimu-  
 57 lus. A linear SVM determines a  $p - 1$  dimensional separating  
 58 hyperplane as a decision boundary in the parameter space, i.e. a  
 59 1-dimensional line for stimuli parameterized in 2-D space,  $\mathbb{R}^2$ .  
 60

The SVM algorithm chooses the optimal separating hyperplane  
 to be the one that maximizes the distance from the hyperplane  
 to the Tone A and Tone B sets.  
 61

Which hyperplane is this? Take the convex hulls of the Tone  
 A and Tone B sets, the set of points enclosed in the tightest  
 rubber band one can stretch around the Tone A and B sets, re-  
 spectively. The optimal hyperplane bisects and is orthogonal to  
 the line segment between the two closest points of the convex  
 hulls (Boyd & Vandenberghe, 2004, p. 46-49). If Tone A and B  
 are linearly inseparable, i.e. if their convex hulls overlap, then a  
 soft margin SVM algorithm can be used, which allows for some  
 points to be on the wrong side of the margin in determining the  
 optimal separating hyperplane, and a soft margin parameter is  
 tuned to balance the tradeoff between maximizing the margin  
 and minimizing classification error.  
 62

We desire the determined classification rule to generalize  
 beyond the particular set of training data used to choose it.  
 Thus, evaluation of classifier performance, e.g. how accurately  
 it identifies tones, is done by determining classification accu-  
 racy on test data, data not in the training data set: in this study,  
 we trained five different classifiers, each one on stimuli from  
 one of the five speakers, and tested the classifiers on the four  
 withheld speakers. For each classifier, we also chose the soft  
 margin parameter to be the value yielding the highest classifi-  
 cation accuracy from a grid search over a set of values ranging  
 from  $1e^{-2}$  to  $1e^2$ .  
 63

## References

- Alexander, J. A. (2010). *The theory of adaptive dispersion and acoustic-phonetic properties of cross-language lexical-tone systems*. Ph.D. thesis Northwestern University.
- Andruski, J. E., & Costello, J. (2004). Using polynomial equations to model pitch contour shape in lexical tones: An example from Green Mong. *Journal of the International Phonetic Association*, 34, 125–140.
- Barnes, J., Veilleux, N., Brugos, A., & Shattuck-Hufnagel, S. (2012). Tonal center of gravity: a global approach to tonal implementation in a level-based intonational phonology. *Laboratory Phonology*, 3, 337–383.
- Barr, D. J., Levy, R., Scheepers, C., & Tily, H. J. (2013). Random effects structure for confirmatory hypothesis testing: keep it maximal. *Journal of Memory and Language*, 68, 255–278.
- Barry, J. G., & Blamey, P. J. (2004). The acoustic analysis of tone differentiation as a means for assessing tone production in speakers of Cantonese. *The Journal of the Acoustical Society of America*, 116, 1739–1748.
- Bashford, J. A., Riener, K. R., & Warren, R. M. (1992). Increasing the intelligibility of speech through multiple phonemic restorations. *Perception and Psychophysics*, 51, 211–217.
- Bashford, J. A., Warren, R. M., & Brown, C. A. (1996). Use of speech-modulated noise adds strong “bottom-up” cues for phonemic restoration. *Perception & Psychophysics*, 58, 342–350.
- Bates, D., Maechler, M., Bolker, B., & Walker, S. (2014). *lme4: Linear mixed-effects models using Eigen and S4*. R package version 1.1-6.
- Bauer, R. S., Kwan-hin, C., & Pak-man, C. (2003). Variation and merger of the rising tones in hong kong cantonese. *Language Variation and Change*, 15, 211–225.
- Becker-Kristal, R. (2010). *Acoustic typology of vowel inventories and Dispersion Theory: insights from a large cross-linguistic corpus*. Ph.D. thesis University of California Los Angeles.
- Bennett, K. P., & Bredensteiner, E. J. (2000). Duality and geometry in SVM classifiers. In *Proceedings of the Seventeenth International Conference on Machine Learning* (pp. 57–64). Morgan Kaufmann Publishers Inc.
- Bishop, J., & Keating, P. (2012). Perception of pitch location within a speaker’s range: Fundamental frequency, voice quality and speaker sex. *Journal of the Acoustical Society of America*, 132, 1100–1112.

- 1 de Boer, B., & Kuhl, P. K. (2003). Investigating the role of infant-directed  
2 speech with a computer model. *Acoustics Research Letters Online*, 4, 129–  
3 134.
- 3 Boersma, P., & Weenink, D. (2010). Praat: doing phonetics by computer (ver-  
4 sion 5.1.32) [computer program]. <http://www.praat.org>.
- 5 Boyd, S., & Vandenberghe, L. (2004). *Convex optimization*. Cambridge Uni-  
6 versity Press.
- 7 Brainard, D. H. (1997). The psychophysics toolbox. *Spatial vision*, 10, 443–  
8 446.
- 9 Burges, C. J. (1998). A tutorial on support vector machines for pattern recog-  
nition. *Data Mining and Knowledge Discovery*, 2, 121–167.
- 10 Chandrasekaran, B., Krishnan, A., & Gandour, J. T. (2007). Experience-  
11 dependent neural plasticity is sensitive to shape of pitch contours. *NeuroReport*, 18, 1963–1967.
- 12 Chang, C.-C., & Lin, C.-J. (2001). LIBSVM: a library  
13 for support vector machines. Software available at  
14 <http://www.csie.ntu.edu.tw/~cjlin/libsvm>.
- 15 Chao, Y.-R. (1930). A system of tone-letters. *Le Maître Phonétique*, 45, 24–27.
- 16 Chao, Y. R. (1968). *A grammar of spoken Chinese*. Berkeley, CA: University  
17 of California Press.
- 18 Chen, M., Yang, Z., & Liu, W. (2014). Deep neural networks for mandarin  
19 tone recognition. In *Neural Networks (IJCNN), 2014 International Joint  
Conference on* (pp. 1154–1158).
- 20 Ciocca, V., & Bregman, A. S. (1987). Perceived continuity of gliding and  
21 steady-state tones through interrupting noise. *Perception and psychophysics*,  
42, 476–484.
- 22 Clarke, C. M., & Garrett, M. F. (2004). Rapid adaptation to foreign-accented  
23 English. *Journal of the Acoustical Society of America*, 116, 3647–3658.
- 24 Cortes, C., & Vapnik, V. (1995). Support-vector networks. *Machine Learning*,  
20, 273–297.
- 25 Dannenbring, G. L. (1976). Perceived auditory continuity with alternately rising  
26 and falling frequency transitions. *Canadian Journal of Psychology*, 30,  
99–114.
- 27 de Leeuw, J., & Mair, P. (2009). Multidimensional scaling using majorization:  
28 SMACOF in R. *Journal of Statistical Software*, 31, 1–30.
- 29 De Looze, C., & Rauzy, S. (2009). Automatic detection and prediction of  
30 topic changes through automatic detection of register variations and pause  
31 duration. In *INTERSPEECH-2009* (pp. 2919–2922).
- 32 DiCanio, C., Amith, J. D., & García, R. C. (2014). The phonetics of moraic  
33 alignment in Yoloxóchitl Mixtec. In *TAL-2014* (pp. 203–210).
- 34 Evanini, K. (2011). Improved measurement point selection for automatic formant  
35 extraction. In *The 85th Annual Meeting of the Linguistic Society of America*. Pittsburgh, PA.
- 36 Evanini, K., & Lai, C. (2010). The importance of optimal parameter setting  
37 for pitch extraction. *The Journal of the Acoustical Society of America*, 128,  
2291.
- 38 Feldman, N. H., Griffiths, T. L., Goldwater, S., & Morgan, J. L. (2013). A role  
39 for the developing lexicon in phonetic category acquisition. *Psychological  
Review*, 120, 751–778.
- 40 Gandour, J. (1981). Perceptual dimensions of tone: evidence from Cantonese.  
41 *Journal of Chinese Linguistics*, 9, 20–36.
- 42 Gandour, J. (1983). Tone perception in Far Eastern languages. *Journal of  
43 Phonetics*, 11, 149–175.
- 44 Gandour, J., Potisuk, S., Dechonkit, S., & Ponglarpisit, S. (1992). Tonal coar-  
45 ticulation in Thai disyllabic utterances: a preliminary study. *Linguistics of the  
46 Tibeto-Burman area*, 15.
- 47 Gandour, J. T., & Harshman, R. A. (1978). Crosslanguage differences in tone  
48 perception: a multidimensional scaling investigation. *Language and Speech*,  
21, 1–33.
- 49 Ganong III, W. F. (1980). Phonetic categorization in auditory word perception.  
50 *Journal of Experimental Psychology: Human Perception and Performance*,  
6, 110–125.
- 51 Gauthier, B., Shi, R., & Xu, Y. (2007). Learning phonetic categories by tracking  
52 movements. *Cognition*, 103, 80–106.
- 53 Gottfried, T. L., & Suiter, T. L. (1997). Effect of linguistic experience on the  
54 identification of mandarin chinese vowels and tones. *Journal of Phonetics*,  
25, 207–231.
- 55 Grassi, M., & Soranzo, A. (2009). MLP: a MATLAB toolbox for rapid and  
56 reliable auditory threshold estimation. *Behavior Research Methods*, 41, 20–  
57 28.
- 58 Greenberg, S., & Zee, E. (1977). On the perception of contour tones. *UCLA  
60*
- 61 Working Papers in Phonetics, 45, 150–159.
- 62 Guyon, I., & Elisseeff, A. (2003). An introduction to variable and feature se-  
63 lection. *Journal of Machine Learning Research*, 3, 1157–1182.
- 64 Han, M., & Kim, K.-O. (1974). Phonetic variation of Vietnamese tones in  
65 disyllabic utterances. *Journal of Phonetics*, 22, 477–492.
- Hant, J. J., & Alwan, A. (2003). A psychoacoustic-masking model to predict  
the perception of speech-like stimuli in noise. *Speech Communication*, 40,  
291–313.
- Hermes, D. J. (2006). Stylization of pitch contours. In S. Sudhoff, D. Lenertová,  
R. Meyer, S. Pappert, P. Augurzky, I. Mleinek, N. Richter, & J. Schließer  
(Eds.), *Methods in empirical prosody research* (pp. 29–62). Walter de  
Gruyter.
- Hess, W. (1983). *Pitch determination of speech signals: algorithms and de-  
vices*. Springer-Verlag.
- Hirst, D., & Espesser, R. (1993). Automatic modelling of fundamental  
frequency using a quadratic spline function. *Travaux de l'Institut de  
Phonétique d'Aix*, 15, 71–85.
- House, D. (1990). *Tonal perception in speech*. Lund, Sweden: Lund University  
Press.
- House, D. (2004a). Pitch and alignment in the perception of tone and intonation.  
In G. Fant, H. Fujisaki, J. Cao, & Y. Xu (Eds.), *From traditional phonology  
to modern speech processing* (pp. 189–204). Foreign Language Teaching  
and Research Press.
- House, D. (2004b). Pitch and alignment in the perception of tone and intonation:  
pragmatic signals and biological codes. In *International Symposium on  
Tonal Aspects of Languages with Emphasis on Tone Languages* (pp. 93–96).
- James, G., Witten, D., Hastie, T., & Tibshirani, R. (2013). *An introduction to  
statistical learning*. New York, NY: Springer.
- Jansen, A., & Niyogi, P. (2009). Point process models for event-based speech  
recognition. *Speech Communication*, (pp. 1155–1168).
- Khouw, E., & Ciocca, V. (2007). Perceptual correlates of cantonese tones.  
*Journal of Phonetics*, 35, 104–117.
- Kochanski, G., Grabe, E., Coleman, J., & Rosner, B. (2005). Loudness pre-  
dicts prominence: Fundamental frequency lends little. *The Journal of the  
Acoustical Society of America*, 118, 1038–1054.
- Krishnan, A., Gandour, J. T., Bidelman, G. M., & Swaminathan, J. (2009).  
Experience-dependent neural representation of dynamic pitch in the brain-  
stem. *NeuroReport*, 20, 408–413.
- Kruschke, J. K. (2013). Bayesian estimation supersedes the t test. *Journal of  
Experimental Psychology: General*, 142, 573–603.
- Kuang, J. (2013). The tonal space of contrastive five level tones. *Phonetica*,  
70, 1–23.
- Ladefoged, P. (1980). What are linguistic sounds made of? *Language*, 56,  
485–502.
- Lam, Y. F. (2014). *Cantonese Tone Recognition Using the Hilbert-Huang  
Transform*. Master's thesis The Hong Kong University of Science and Tech-  
nology Hong Kong, China.
- Laniran, Y. O. (1992). *Intonation in tone languages: the phonetic implementa-  
tion of tones in Yoruba*. Ph.D. thesis Cornell University.
- Lee, C. (2009). Identifying isolated, multispeaker mandarin tones from brief  
acoustic input: A perceptual and acoustic study. *The Journal of the Acousti-  
cal Society of America*, 125, 1125–1137.
- Lee, C., Tao, L., & Bond, Z. (2008). Identification of acoustically modified  
mandarin tones by native listeners. *Journal of Phonetics*, 36, 537–563.
- Lee, C., Tao, L., & Bond, Z. (2009). Speaker variability and context in the iden-  
tification of fragmented mandarin tones by native and non-native listeners.  
*Journal of Phonetics*, 37, 1–15.
- Levow, G.-A. (2005). Context in multi-lingual tone and pitch accent recogni-  
tion. In *Proceedings of INTERSPEECH 2005* (pp. 1809–1812).
- Levow, G.-A. (2006). Unsupervised and semi-supervised learning of tone and  
pitch accent. In *Proceedings of the Human Language Technology Confer-  
ence of the North American Chapter of the ACL* (pp. 224–231).
- Li, Q., & Chen, Y. (2016). An acoustic study of contextual tonal variation in  
Tianjin Mandarin. *Journal of Phonetics*, 54, 123 – 150.
- Li, Y., & Lee, T. (2007). Perceptual equivalence of approximated Cantonese  
tone contours. In *INTERSPEECH-2007* (pp. 2677–2680).
- Li, Y., & Lee, T. (2008). A perceptual study of approximated cantonese tone  
contours. In *Chinese Spoken Language Processing, 2008. ISCSLP '08. 6th  
International Symposium on* (pp. 1–4).
- Li, Y., Lee, T., & Qian, Y. (2002). Acoustical F0 analysis of continuous Can-  
tonese speech. In *Proceedings of International Symposium on Chinese Spo-*

- ken Language Processing (pp. 127–130).
- Li, Y., Lee, T., & Qian, Y. (2004). F0 analysis and modeling for Cantonese Text-to-Speech. In *SP-2004* (pp. 467–470).
- Liljencrants, J., & Lindblom, B. (1972). Numerical simulation of vowel quality systems: The role of perceptual contrast. *Language*, 48, 839–862.
- Lisker, L. (1978). Rapid vs. rabid: A catalogue of acoustic features that may cue the distinction. *Haskins Laboratories Status Report on Speech Research*, SR-54, 127–132.
- Liu, S., & Samuel, A. G. (2004). Perception of mandarin lexical tones when f0 information is neutralized. *Language and Speech*, 47, 109–138.
- Luce, R. D. (1963). Detection and recognition. In R. Luce, R. R. Bush, & E. Galanter (Eds.), *Handbook of mathematical psychology* (pp. 103–189). Wiley volume 1.
- Matthews, S., & Yip, V. (1994). *Cantonese: a comprehensive grammar*. New York: Routledge.
- McGurk, H., & MacDonald, J. (1976). Hearing lips and seeing voices. *Nature*, 264, 746–748.
- Meyer, D., Dimitriadou, E., Hornik, K., Weingessel, A., & Leisch, F. (2012). *e1071: Misc Functions of the Department of Statistics (e1071)*. TU Wien. R package version 1.6-1.
- Miller, G. A., & Licklider, J. C. R. (1950). The intelligibility of interrupted speech. *The Journal of the Acoustical Society of America*, 22, 167–173.
- Mok, P. P.-K., & Wong, P. W.-Y. (2010a). Perception of the merging tones in Hong Kong Cantonese: preliminary data on monosyllables. In *INTERSPEECH-2010*.
- Mok, P. P.-K., & Wong, P. W.-Y. (2010b). Production of the merging tones in Hong Kong Cantonese: preliminary data on monosyllables. In *INTERSPEECH-2010*.
- Morén, B., & Zsiga, E. (2006). The lexical and Post-Lexical phonology of thai tones\*. *Natural Language & Linguistic Theory*, 24, 113–178.
- Nearey, T. M., & Assmann, P. F. (1986). Modeling the role of inherent spectral change in vowel identification. *The Journal of the Acoustical Society of America*, 80, 1297–1308.
- Nosofsky, R. M. (1990). Relations between exemplar-similarity and likelihood models of classification. *Journal of Mathematical Psychology*, 34, 393–418.
- Pelli, D. G. (1997). The videotoolbox software for visual psychophysics: transforming numbers into movies. *Spatial vision*, 10, 437–442.
- Peng, G., & Wang, W. S. (2005). Tone recognition of continuous cantonese speech based on support vector machines. *Speech Communication*, 45, 49–62.
- Peng, G., Zheng, H., & Wang, W. S. Y. (2004). Tone recognition for chinese speech: a comparative study of mandarin and cantonese. In *Chinese Spoken Language Processing, 2004 International Symposium on* (pp. 233–236).
- Peterson, G. E., & Barney, H. L. (1952). Control methods used in a study of the vowels. *The Journal of the Acoustical Society of America*, 24, 175–184.
- Pisarn, C., & Theeramunkong, T. (2006). Improving Thai spelling recognition with tone features. In T. Salakoski, F. Ginter, S. Pyysalo, & T. Pahikkala (Eds.), *Advances in Natural Language Processing: 5th International Conference on NLP, FinTAL 2006 Turku, Finland, August 23–25, 2006 Proceedings* (pp. 388–398). Berlin, Heidelberg: Springer Berlin Heidelberg.
- Prom-on, S., Xu, Y., & Thipakorn, B. (2009). Modeling tone and intonation in Mandarin and English as a process of target approximation. *Journal of the Acoustical Society of America*, 125, 405–424.
- Prukkanon, N., Channongthai, K., & Miyanaga, Y. (2016). {F0} contour approximation model for a one-stream tonal word recognition system. *{AEU} - International Journal of Electronics and Communications*, 70, 681 – 688.
- Qian, Y., Lee, T., & Soong, F. K. (2007). Tone recognition in continuous cantonese speech using supratone models. *The Journal of the Acoustical Society of America*, 121, 2936–2945.
- R Core Team (2014). *R: A Language and Environment for Statistical Computing*. R Foundation for Statistical Computing Vienna, Austria.
- Remijzen, B. (2013). Tonal alignment is contrastive in falling contours in Dinka. *Language*, 89, 297–327.
- Remijzen, B., & Ayoker, O. G. (2014). Contrastive tonal alignment in falling contours in shilluk. *Phonology*, 31, 435–462.
- Rosenfelder, I., Fruehwald, J., Evanini, K., & Yuan, J. (2011). FAVE (forced alignment and vowel extraction) program suite. <http://fave.ling.upenn.edu>.
- Samuel, A. (1996). Phoneme restoration. *Language and Cognitive Processes*, 11, 647.
- Shen, J., Deutsch, D., & Rayner, K. (2013). On-line perception of Mandarin Tones 2 and 3: evidence from eye movements. *Journal of the Acoustical Society of America*, 133, 3016–3029.
- Shih, C., & Lu, H.-Y. D. (2015). Effects of talker-to-listener distance on tone. *Journal of Phonetics*, 51, 6 – 35. What's So Special About H(igh)? Multi-Disciplinary Perspectives on the Linguistic Functions of Raised Pitch.
- Silbert, N. (2014). Visualizing confusion matrices. <Http://www.nhsilbert.net/source/2014/03/visualizing-confusion-matrices/>.
- Strange, W., Jenkins, J. J., & Johnson, T. L. (1983). Dynamic specification of coarticulated vowels. *The Journal of the Acoustical Society of America*, 74, 695–705.
- Sumner, M., & Samuel, A. G. (2009). The effect of experience on the perception and representation of dialect variants. *Journal of Memory and Language*, 60, 487–501.
- Talkin, D. (1995). A robust algorithm for pitch tracking (RAPT). In W. B. Kleijn, & K. K. Paliwal (Eds.), *Speech coding and synthesis* (pp. 495–518). Elsevier Science Inc.
- Taylor, P. (2000). Analysis and synthesis of intonation using the tilt model. *The Journal of the Acoustical Society of America*, 107, 1697–1714.
- Tian, Y., Zhou, J., Chu, M., & Chang, E. (2004). Tone recognition with fractionized models and outlined features. In *Acoustics, Speech, and Signal Processing, 2004. Proceedings. (ICASSP '04). IEEE International Conference on* (pp. I–105–8 vol.1). volume 1.
- Vallabha, G. K., McClelland, J. L., Pons, F., Werker, J. F., & Amano, S. (2007). Unsupervised learning of vowel categories from infant-directed speech. *Proceedings of the National Academy of Sciences*, 104, 13273–13278.
- Vapnik, V. N. (1995). *The nature of statistical learning*. Springer.
- Wang, S., & Levow, G.-A. (2008). Mandarin Chinese tone nucleus detection with landmarks. In *Proceedings of Interspeech 2008* (pp. 1101–1104).
- Wang, S., Tang, Z., Zhao, Y., & Ji, S. (2009). Tone recognition of continuous Mandarin speech based on binary-class SVMs. In *Information Science and Engineering (ICISE), 2009 1st International Conference on* (pp. 710–713).
- Warren, R. M. (1970). Perceptual restoration of missing speech sounds. *Science*, 167, 392–393.
- Whalen, D. H., & Xu, Y. (1992). Information for Mandarin tones in the amplitude contour and in brief segments. *Phonetica*, 49, 25–47.
- Wickham, H. (2009). *ggplot2: elegant graphics for data analysis*. Springer.
- Wong, P. C. M., & Diehl, R. L. (2003). Perceptual normalization for inter- and intratalker variation in cantonese level tones. *Journal of Speech, Language & Hearing Research*, 46, 413–421.
- Wong, Y. W. (2006a). Contextual tonal variations and pitch targets in Cantonese. In *Proceedings of Speech Prosody 2006, Dresden*.
- Wong, Y. W. (2006b). Realization of cantonese rising tones under different speaking rates. In *Proceedings of Speech Prosody 2006, Dresden*.
- Xu, Y. (1997). Contextual tonal variations in Mandarin. *Journal of Phonetics*, 25, 61–83.
- Xu, Y. (1998). Consistency of tone-syllable alignment across different syllable structures and speaking rates. *Phonetica*, 55, 179–203.
- Xu, Y. (2001). Fundamental frequency peak delay in Mandarin. *Phonetica*, 58, 26–52.
- Yu, K. M., & Lam, H. W. (2011). The role of creaky voice in Cantonese tonal perception. In *Proceedings of ICPhS XVII*.
- Yu, K. M., & Lam, H. W. (2014). The role of creaky voice in Cantonese tonal perception. *Journal of the Acoustical Society of America*, 136, 1320–1333.
- Zhang, J., & Hirose, K. (2000). Anchoring hypothesis and its application to tone recognition of chinese continuous speech. In *Acoustics, Speech, and Signal Processing, 2000. ICASSP '00. Proceedings. 2000 IEEE International Conference on* (pp. 1419–1422). volume 3.
- Zhang, J., & Hirose, K. (2004). Tone nucleus modeling for Chinese lexical tone recognition. *Speech Communication*, 42, 447–466.
- Zhou, N., Zhang, W., Lee, C., & Xu, L. (2008). Lexical tone recognition with an artificial neural network. *Ear and hearing*, 29, 326–335. PMC2562432.
- Zsiga, E., & Nitisoroj, R. (2007). Tone features, tone perception, and peak alignment in thai. *Language and Speech*, 50, 343–383.

**Supplementary Material**

[\*\*Click here to download Supplementary Material: res-jphon-supp-clean.pdf\*\*](#)