

Representation Learning for Discovering Phonemic Tone Contours

Bai Li^{1,2}, Jing Yi Xie¹, Frank Rudzicz^{1,2,3}

¹ University of Toronto, Toronto, Canada

² Vector Institute, Toronto, Canada

³ St Michael's Hospital, Toronto, Canada

{bai, frank}@cs.toronto.edu, jingyi.xie@mail.utoronto.ca

Abstract

Tone is a prosodic feature used to distinguish words in many languages, some of which are endangered and scarcely documented. In this work, we use unsupervised representation learning to identify probable clusters of syllables that share the same phonemic tone. Our method extracts the pitch for each syllable, then trains a convolutional autoencoder to learn a low-dimensional representation for each contour. We then apply the mean shift algorithm to cluster tones in high-density regions of the latent space. Furthermore, by feeding the centers of each cluster into the decoder, we produce a prototypical contour that represents each cluster. We apply this method to spoken multi-syllable words in Mandarin Chinese and Cantonese and evaluate how closely our clusters match the ground truth tone categories. Finally, we discuss some difficulties with our approach, including contextual tone variation and allophony effects.

1 Introduction

Tonal languages use pitch to distinguish different words, for example, *yi* in Mandarin may mean ‘one’, ‘to move’, ‘already’, or ‘art’, depending on the pitch contour. Of over 6000 languages in the world, it is estimated that as many as 60-70% are tonal (Lewis, 2009; Yip, 2002). A few of these are national languages (e.g., Mandarin Chinese, Vietnamese, and Thai), but many tonal languages have a small number of speakers and are scarcely documented. There is a limited availability of trained linguists to perform language documentation before these languages become extinct, hence the need for better tools to assist linguists in these tasks.

One of the first tasks during the description of an unfamiliar language is determining its phonemic inventory: what are the consonants, vowels, and tones of the language, and which pairs of phonemes

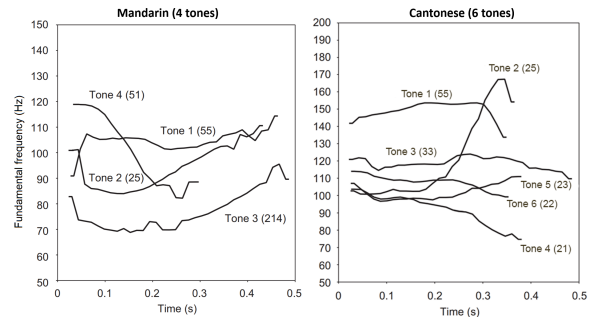


Figure 1: Fundamental frequency (F0) contours for the four Mandarin tones and six Cantonese tones in isolation, produced by native speakers. Figure adapted from (Francis et al., 2008).

are contrastive? Tone presents a unique challenge because unlike consonants and vowels, which can be identified in isolation, tones do not have a fixed pitch, and vary by speaker and situation. Since tone data is subject to interpretation, different linguists may produce different descriptions of the tone system of the same language (Yip, 2002).

In this work, we present a model to automatically infer phonemic tone categories of a tonal language. We use an unsupervised learning approach: a convolutional autoencoder learns a low-dimensional representation of each tone using only a set of spoken syllables in the target language. This is followed by mean shift clustering to identify clusters of syllables that probably have the same tone. We apply our method on Mandarin Chinese and Cantonese datasets, for which the ground truth annotation is used for evaluation. Our method does not make any language-specific assumptions, so it may be applied to low-resource languages whose phonemic inventories are not already established.

1.1 Tone in Mandarin and Cantonese

Mandarin Chinese (1.1 billion speakers) and Cantonese (74 million speakers) are two tonal lan-

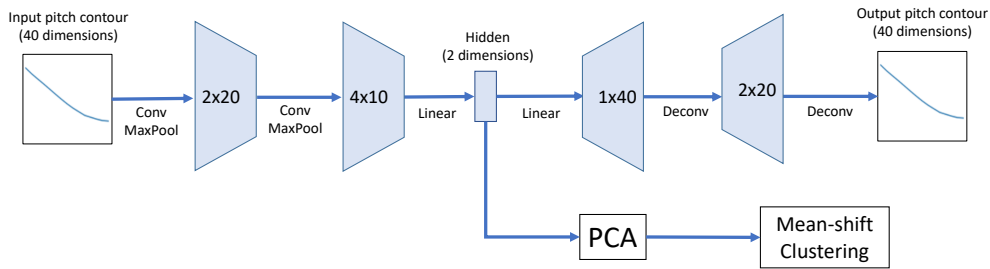


Figure 2: Diagram of our model architecture, consisting of a convolutional autoencoder to learn a latent representation for each pitch contour, and mean shift clustering to identify groups of similar tones.

guages in the Sinitic family (Lewis, 2009). Mandarin has four lexical tones: high (55), rising (25), low-dipping (214), and falling (51)¹. The third tone sometimes undergoes sandhi, addressed in section 3. We exclude a fifth, neutral tone, which can only occur in word-final positions and has no fixed pitch.

Cantonese has six lexical tones: high-level (55), mid-rising (25), mid-level (33), low-falling (21), low-rising (23), and low-level (22). Some descriptions of Cantonese include nine tones, of which three are *checked* tones that are flat, shorter in duration, and only occur on syllables ending in /p/, /t/, or /k/. Since each one of the checked tones are in complementary distribution with an unchecked tone, we adopt the simpler six tone model that treats the checked tones as variants of the high, mid, and low level tones. Contours for the lexical tones in both languages are shown in Figure 1.

2 Related work

Many low-resource languages lack sufficient transcribed data for supervised speech processing, thus unsupervised models for speech processing is an emerging area of research. The Zerospeech 2015 and 2017 challenges featured unsupervised learning of contrasting phonemes in English and Xitsonga, evaluated by an ABX phoneme discrimination task (Versteegh et al., 2015). One successful approach used denoising and correspondence autoencoders to learn a representation that avoided capturing noise and irrelevant inter-speaker variation (Renshaw et al., 2015). Deep LSTMs for segmenting and clustering phonemes in speech have also been explored in (Müller et al., 2017b) and (Müller et al., 2017a).

In Mandarin Chinese, deep neural networks have been successful for tone classification in isolated

¹The numbers are Chao tone numerals, where 1 is the lowest and 5 is the highest pitch.

syllables (Chen et al., 2016) as well as in continuous speech (Ryant et al., 2014b,a). Both of these models found that Mel-frequency cepstral coefficients (MFCCs) outperformed pitch contour features, despite the fact that MFCC features do not contain pitch information. In Cantonese, support vector machines (SVMs) have been applied to classify tones in continuous speech, using pitch contours as input (Peng and Wang, 2005).

Unsupervised learning of tones remains largely unexplored. Levow (2006) performed unsupervised and semi-supervised tone clustering in Mandarin, using average pitch and slope as features, and k -means and asymmetric k -lines for clustering. Graph-based community detection techniques have been applied to group n -grams of contiguous contours into clusters in Mandarin (Zhang, 2019). In recent work concurrent to ours, Fry (2020) uses adversarial autoencoders and hierarchical clustering to identify tone inventories, and evaluate their method on Mandarin, Cantonese, Fungwa, and English data.

We further explore unsupervised deep neural networks for phonemic tone clustering. It should be noted that our unsupervised model is not given tone labels during training, and the number of tones is assumed to be unknown, so it cannot be directly compared to supervised tone classifiers in the literature.

3 Data and preprocessing

We use data from Mandarin Chinese and Cantonese. For each language, the data consists of a list of spoken words, recorded by the same speaker. The Mandarin dataset is from a female speaker and is provided by Shtooka², and the Cantonese dataset is from a male speaker and is downloaded from

²<http://shtooka.net/>, specifically the cmn-caentan dataset.

Forvo³, an online crowd-sourced pronunciation dictionary. We require all samples within each language to be from the same speaker to avoid the difficulties associated with channel effects and inter-speaker variation. We randomly sample 400 words from each language, which are mostly between 2 and 4 syllables; to reduce the prosody effects with longer utterances, we exclude words longer than 4 syllables.

We extract ground-truth tones for evaluation purposes. In Mandarin, the tones are extracted from the pinyin transcription; in Cantonese, we reference the character entries on Wiktionary⁴ to retrieve the romanized pronunciation and tones. For Mandarin, we adjust for third-tone sandhi (a phonological rule where a pair of consecutive third-tones is always realized as a second-tone followed by a third-tone), and use the sandhi tone as the ground truth. We also exclude the neutral tone, which has no fixed pitch and is sometimes thought of as a lack of tone.

3.1 Pitch extraction and syllable segmentation

We use Praat’s autocorrelation-based pitch estimation algorithm to extract the fundamental frequency (F0) contour for each sample, using a minimum frequency of 75Hz and a maximum frequency of 500Hz (Boersma, 1993). The interface between Python and Praat is handled using Parselmouth (Jadoul et al., 2018). We normalize the contour to be between 0 and 1, based on the speaker’s pitch range.

Next, we manually segment each speech sample into syllables, necessary because syllable boundaries are not provided in our datasets. We sample the pitch at 40 equally spaced points, obtaining a constant length vector as input to our model. Note that by sampling a variable length contour to a constant length, the model does not have information about syllable length; we discuss this design choice in section 6.2.

4 Model

4.1 Convolutional autoencoder

We use a convolutional autoencoder (Figure 2) to learn a two-dimensional latent vector for each syllable. Convolutional layers are widely used in computer vision and speech processing to learn spatially local features that are invariant of position.

³<https://forvo.com/>

⁴<https://en.wiktionary.org/>

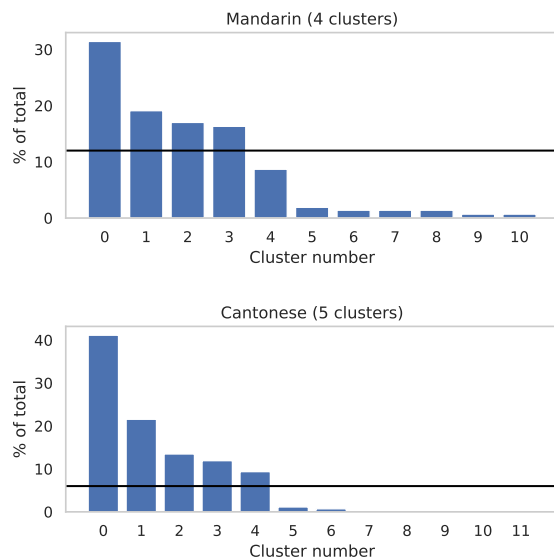


Figure 3: Clusters generated by the mean shift procedure. The black line shows the threshold: we discard clusters with size below this value and treat their points as unclustered.

We use a low dimensional latent space so that the model learns to generate a representation that only captures the most important aspects of the input contour, and also because clustering algorithms tend to perform poorly in high dimensional spaces.

Our encoder consists of three layers. The first layer applies 2 convolutional filters (kernel size 4, stride 1) followed by max pooling (kernel size 2) and a tanh activation. The second layer applies 4 convolutional filters (kernel size 4, stride 1), again with max pooling (kernel size 2) and a tanh activation. The third layer is a fully connected layer with two dimensional output. Our decoder is the encoder in reverse, consisting of one fully connected layer and two deconvolution layers, with the same layer shapes as the encoder.

We train the autoencoder using PyTorch (Paszke et al., 2017), for 500 epochs, with a batch size of 60. The model is optimized using Adam (Kingma and Ba, 2015) with a learning rate of 5e-4 to minimize the mean squared error between the input and output contours.

4.2 Mean shift clustering

We run the encoder on each syllable’s pitch contour to get their latent representations; we apply principal component analysis (PCA) to remove any correlation between the two dimensions. Then, we run mean shift clustering (Comaniciu and Meer, 2002; Ghassabeh and Rudzicz, 2018), estimating

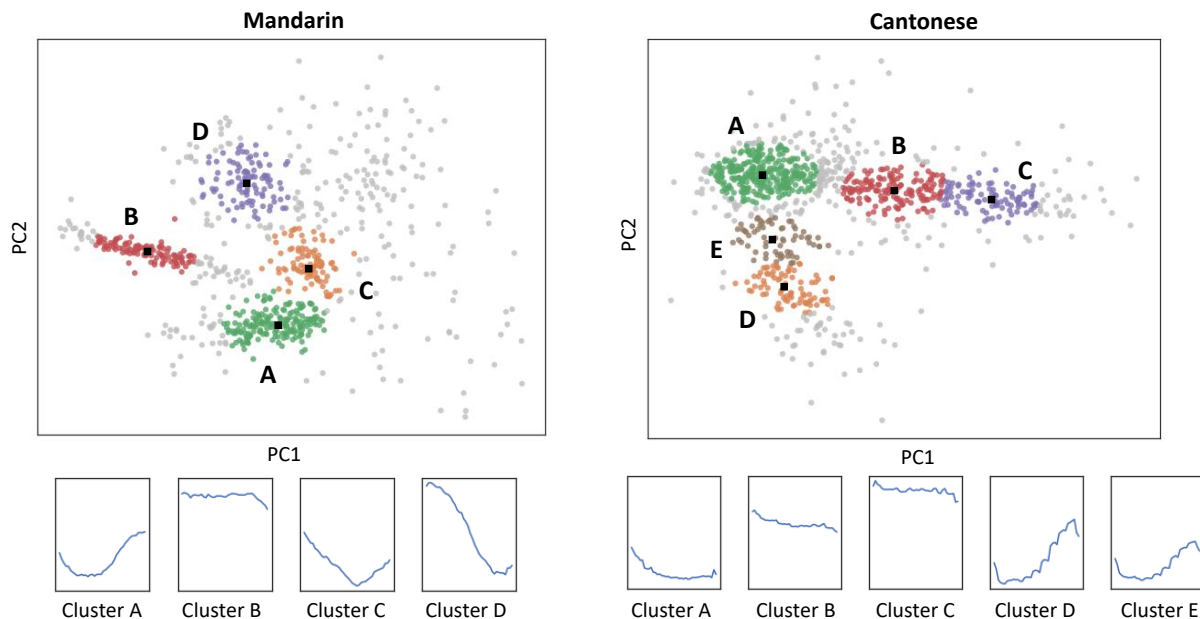


Figure 4: Latent space generated by autoencoder and the results of mean shift clustering for Mandarin and Cantonese. Each cluster center is fed through the decoder to generate the corresponding pitch contour. The clusters within each language are ordered by size, from largest to smallest.

a probability density function in the latent space. The procedure performs gradient ascent on all the points until they converge to a set of stationary points, which are local maxima of the density function. These stationary points are taken to be cluster centers, and points that converge to the same stationary point belong to the same cluster. We feed the cluster centers into the decoder to generate a prototype pitch contour for each cluster.

Unlike k -means clustering, the mean shift procedure does not require the number of clusters to be specified, only a bandwidth parameter (set to 0.6 for our experiments). The cluster centers are always in regions of high density, so they can be viewed as prototypes that represent their respective clusters. Another advantage is that unlike k -means, mean shift clustering is robust to outliers.

4.3 Selecting bandwidth and threshold

The bandwidth parameter controls the size of the clusters: a higher bandwidth value generates fewer and larger clusters. We tune the bandwidth parameter to produce linguistically plausible tone clusters: we expect between 3 to 8 different clusters, each clusters should have at least 1/10 of the points be assigned to it, and most points should belong to some cluster.

The mean shift procedure assigns every point to some cluster, even if the resulting cluster contains

only a few points. Thus, we set a threshold: we treat clusters smaller than the threshold as spurious, and leave their points as unclustered. Figure 3 shows the effect of the threshold on both languages.

4.4 k -means baseline

We implement a simple k -means baseline similar to Levow (2006), using two engineered features. The first feature is the average pitch of all the points in the pitch contour; the second feature is the slope of an ordinary least squares regression fit on the pitch contour. After extracting these features for every syllable, we run k -means clustering, using the same number of clusters that is chosen by the mean shift algorithm.

5 Results

Figure 4 shows the latent space learned by the autoencoders and the clustering output. Our model found 4 tone clusters in Mandarin, matching the number of phonemic tones (Table 1) and 5 in Cantonese, which is one fewer than the number of phonemic tones (Table 2). In Mandarin, the 4 clusters correspond very well with the the 4 phonemic tone categories, and the generated contours closely match the ground truth in Figure 1. There is some overlap between tones 3 and 4; this is because tone 3 is sometimes realized a low-falling tone without the final rise, a process known as half T3 sandhi

Cluster	T1	T2	T3	T4
A	1	163	12	4
B	108	0	0	1
C	0	5	53	31
D	1	0	0	97
N/A	47	30	53	129

Table 1: Cluster and tone frequencies for Mandarin.

Cluster	T1	T2	T3	T4	T5	T6
A	5	5	59	109	7	105
B	102	3	36	2	2	7
C	93	0	0	2	0	0
D	0	64	4	3	2	11
E	0	28	2	4	30	2
N/A	70	39	51	45	15	49

Table 2: Cluster and tone frequencies for Cantonese.

(Chen, 2000), thus, it may overlap with tone 4 (falling tone).

In Cantonese, the 5 clusters A-E correspond to low-falling, mid-level, high-level, mid-rising, and low-rising tones. Tone clustering in Cantonese is expected to be more difficult than in Mandarin because of 6 contrastive tones, rather than 4. The model is more effective at clustering the higher tones (1, 2, 3), and less effective at clustering the lower tones (4, 5, 6), particularly tone 4 (low-falling) and tone 6 (low-level). This confirms the difficulties in prior work, which reported worse classification accuracy on the lower-pitched tones because the lower region of the Cantonese tone space is more crowded than the upper region (Peng and Wang, 2005).

To evaluate how much the clusters match the ground truth, we use normalized mutual information (NMI); this is preferable over accuracy because it does not require the number of detected clusters to be the same as the number of tones. In Table 3, we evaluate NMI for our autoencoder model and the k -means baseline. We consider two scenarios for each language: using all the syllables (All) and using only the first syllable of each word (First).

In all cases, the clusters from the autoencoder model have higher NMI than the k -means model. The improvement is due to the mean shift procedure identifying points that belong to a cluster with high confidence: it only makes predictions for those points, whereas k -means assigns every point to a cluster. All models perform better on the

	Autoencoder	k -means
Mandarin (First)	0.846	0.829
Mandarin (All)	0.753	0.645
Cantonese (First)	0.575	0.493
Cantonese (All)	0.463	0.377

Table 3: Normalized mutual information (NMI) between cluster assignments and ground truth tones, considering only the first syllable of each word, or all syllables.

first syllable of each utterance than the rest of the syllables; we discuss the reasons for this in the next section.

6 Limitations

6.1 Contextual effects

One limitation of our model is it considers syllables in isolation, but in reality, pitch is affected by context. Two types of contextual effects are carry-over and declination. A carry-over effect is when the pitch contour of a tone undergoes contextual variation depending on the preceding tone; strong carry-over effects have been observed in Mandarin (Xu, 1997). Prior work (Levow, 2006) avoided carry-over effects by using only the second half of every syllable, but we do not consider language-specific heuristics in our model.

Declination is a phenomenon in which the pitch declines over an utterance (Yip, 2002; Peng and Wang, 2005). This is especially a problem in Cantonese, which has tones that differ only on pitch level and not contour: for example, a mid-level tone near the end of a phrase may have the same absolute pitch as a low-level tone at the start of a phrase.

Contextual effects are apparent in our results (Table 3). In both Mandarin and Cantonese, the clustering is more accurate when using only the first syllable (which is not affected by carry-over or declination), compared to using all the syllables.

6.2 Minimal pairs and allotones

Tone is not a purely phonetic property: it is impossible to determine, from phonetics alone, whether two pitch contours have the same or different tones. The same underlying tone may manifest as several different allotones depending on the phonetic context.

An example of this appears in Cantonese. Its tone system is sometimes analyzed as having nine

tones instead of six, where six of the tones are only permitted in open syllables (e.g. *si*) and three are only permitted in checked syllables (e.g. *sik*). Other analyses use a six-tone system, treating the three checked tones as allotonic variants of the high, mid, and low tones. By taking this approach, one implies that length is a property of the syllable and cannot be solely responsible for contrasting two tones.

Length is not the only differentiating factor for allotones. Another example is in Wu Chinese, where syllables beginning with voiced consonants have lower pitch than those beginning with voiceless consonants (Yip, 2002). Thus the same language may have vastly different numbers of tones, depending on the analysis.

Linguistically, two phonemic tones are considered to be contrastive if there exists a minimal pair: two semantically different lexical items that are identical in every aspect except for tone. This definition is the most widely used because it clearly settles disagreements about whether two tones are same or different. However, it is problematic for unsupervised models that only have access to phonetic and not semantic information. This issue is not unique to tone: similar difficulties have been noted when attempting to identify consonant and vowel phonemes automatically (Kempton and Moore, 2014).

7 Conclusion

We propose a model for unsupervised clustering and discovery of phonemic tones in tonal languages, using spoken words as input. Our model extracts the F0 pitch contour, trains a convolutional autoencoder to learn a low-dimensional representation for each contour, and applies mean shift clustering to the resulting latent space. We obtain promising results with both Mandarin Chinese and Cantonese, using only 400 spoken words from each language. Cantonese presents more difficulties because of its larger number of tones, especially at the lower half of the pitch range, and also due to multiple contrastive level tones. Still, in both our languages, our method finds clusters of tones that better match the ground truth than the *k*-means baseline. Finally, we discuss the effects of contextual variation and the limitations of unsupervised learning for the tone induction problem.

8 Acknowledgments

We thank Prof Gerald Penn for his help suggestions during this project. Rudzicz is a CIFAR Chair in AI.

References

- Paul Boersma. 1993. Accurate short-term analysis of the fundamental frequency and the harmonics-to-noise ratio of a sampled sound. In *Proceedings of the institute of phonetic sciences*, volume 17, pages 97–110. Amsterdam.
- Charles Chen, Razvan C Bunescu, Li Xu, and Chang Liu. 2016. Tone classification in Mandarin Chinese using convolutional neural networks. In *INTER-SPEECH*, pages 2150–2154.
- Matthew Y Chen. 2000. *Tone sandhi: Patterns across Chinese dialects*, volume 92. Cambridge University Press.
- Dorin Comaniciu and Peter Meer. 2002. Mean shift: A robust approach toward feature space analysis. *IEEE Transactions on Pattern Analysis & Machine Intelligence*, (5):603–619.
- Alexander L Francis, Valter Ciocca, Lian Ma, and Kimberly Fenn. 2008. Perceptual learning of Cantonese lexical tones by tone and non-tone language speakers. *Journal of Phonetics*, 36(2):268–294.
- Michael David Fry. 2020. *Grammaticus ex machina: tone inventories as hypothesized by machine*. Ph.D. thesis, University of British Columbia.
- Y Aliyari Ghassabeh and F Rudzicz. 2018. Modified mean shift algorithm. *IET Image Processing*, 12(12):2172–2177.
- Yannick Jadoul, Bill Thompson, and Bart De Boer. 2018. Introducing Parselmouth: A Python interface to Praat. *Journal of Phonetics*, 71:1–15.
- Timothy Kempton and Roger K Moore. 2014. Discovering the phoneme inventory of an unwritten language: A machine-assisted approach. *Speech Communication*, 56:152–166.
- Diederik P Kingma and Jimmy Ba. 2015. Adam: A method for stochastic optimization. *ICLR*.
- Gina-Anne Levow. 2006. Unsupervised and semi-supervised learning of tone and pitch accent. In *Proceedings of the main conference on Human Language Technology Conference of the North American Chapter of the Association of Computational Linguistics*, pages 224–231. Association for Computational Linguistics.
- M. Paul Lewis. 2009. *Ethnologue: Languages of the World*, 16th edition. SIL International, Dallas, Texas.

- Markus Müller, Jörg Franke, Sebastian Stüker, and Alex Waibel. 2017a. Improving phoneme set discovery for documenting unwritten languages. *Elektronische Sprachsignalverarbeitung (ESSV)*, 2017.
- Markus Müller, Jörg Franke, Alex Waibel, and Sebastian Stüker. 2017b. Towards phoneme inventory discovery for documentation of unwritten languages. In *2017 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 5200–5204. IEEE.
- Adam Paszke, Sam Gross, Soumith Chintala, Gregory Chanan, Edward Yang, Zachary DeVito, Zeming Lin, Alban Desmaison, Luca Antiga, and Adam Lerer. 2017. Automatic differentiation in PyTorch. In *NIPS Autodiff Workshop*.
- Gang Peng and William S-Y Wang. 2005. Tone recognition of continuous Cantonese speech based on support vector machines. *Speech Communication*, 45(1):49–62.
- Daniel Renshaw, Herman Kamper, Aren Jansen, and Sharon Goldwater. 2015. A comparison of neural network methods for unsupervised representation learning on the zero resource speech challenge. In *Sixteenth Annual Conference of the International Speech Communication Association*.
- Neville Ryant, Malcolm Slaney, Mark Liberman, Elizabeth Shriberg, and Jiahong Yuan. 2014a. Highly accurate Mandarin tone classification in the absence of pitch information. In *Proceedings of Speech Prosody*, volume 7.
- Neville Ryant, Jiahong Yuan, and Mark Liberman. 2014b. Mandarin tone classification without pitch tracking. In *2014 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 4868–4872. IEEE.
- Maarten Versteegh, Roland Thiollie, Thomas Schatz, Xuan Nga Cao, Xavier Anguera, Aren Jansen, and Emmanuel Dupoux. 2015. The zero resource speech challenge 2015. In *Sixteenth Annual Conference of the International Speech Communication Association*.
- Yi Xu. 1997. Contextual tonal variations in Mandarin. *Journal of phonetics*, 25(1):61–83.
- Moira Yip. 2002. *Tone*. Cambridge University Press.
- Shuo Zhang. 2019. Data mining Mandarin tone contour shapes. *SIGMORPHON 2019*, page 144.