

# Encoding of lexical tone in self-supervised models of spoken language

Anonymous ACL submission

## Abstract

Interpretability research has shown that self-supervised Spoken Language Models (SLMs) encode a wide variety of features in human speech from the acoustic, phonetic, phonological, syntactic and semantic levels, to speaker characteristics. The bulk of prior research on representations of phonology has focused on segmental features such as phonemes; the encoding of suprasegmental phonology (such as tone and stress patterns) in SLMs is not yet well understood. Tone is a suprasegmental feature that is present in more than half of the world’s languages. This paper aims to analyze the tone encoding capabilities of SLMs, using Mandarin and Vietnamese as case studies. We show that SLMs encode lexical tone to a significant degree even when they are trained on data from non-tonal languages. We further find that SLMs behave similarly to native and non-native human participants in tone and consonant perception studies, but they do not follow the same developmental trajectory.

## 1 Introduction

Explaining the inner workings of self-supervised models of written and spoken language has been the focus of much recent work. Transformer-based (Vaswani et al., 2017) written language models have been shown to encode many types of linguistic information (Conneau et al., 2018; Hewitt and Manning, 2019). The analysis of self-supervised Spoken Language Models (SLMs) is also gaining traction: architectures such as wav2vec2 (Baevski et al., 2020) and HuBERT (Hsu et al., 2021) have been shown to encode linguistic information at the phonetic, phonological, syntactic and semantic levels of human speech without labeled data (Abdullah et al., 2021; Ma et al., 2021; de Seyssel et al., 2022; Bartelds et al., 2022; Martin et al., 2023; Shen et al., 2023).

The majority of research on representations of phonetic and phonological information in SLMs

focuses on the segmental level. **Segmental** refers to units of speech that do not spread but remain localized. Phonemes (e.g. vowels and consonants) are the smallest abstract units of sound that help to distinguish one unit from another (e.g. **pat** vs **bat**). **Suprasegmental**, in contrast, refers to features that are not necessarily limited to single units, but can spread across multiple phonemes or phrases. Examples include tone, stress patterns, and intonation, which can all entail syllable and phrase level changes (Singh and Fu, 2016). The representation of suprasegmental information in SLMs is important to study, as it is one of the main distinguishing features of speech compared to text: spoken utterances use suprasegmental cues to convey information that is generally not explicitly marked in a corresponding written sentence. As a first step, in this work, we focus on lexical tone as a highly constrained, relatively well-understood example of a suprasegmental feature.

We firstly examine to what extent SLMs trained on tonal and non-tonal languages encode tone information in their internal representations. We find that SLMs are capable of capturing tonal information, regardless of whether they are trained on tonal or non-tonal languages.

Secondly, we investigate the impact of supervised fine-tuning on the automatic speech recognition (ASR) task. We find that fine-tuning *enhances* tone representations for models trained on tonal languages, but *reduces* them for models trained on non-tonal languages.

Thirdly, we investigate whether SLMs exhibit the same perceptual patterns as native and non-native human listeners. We find that models show patterns similar to humans in discrimination of Mandarin tones and consonants, but find no evidence that they follow a similar developmental trajectory.

## 2 Tones

Estimates suggest that more than 60% of the world’s languages use some degree of tonal contrast (Yip, 2002). Our primary focus is on lexical tone, the process by which lexical items are distinguished from one another primarily by pitch cues (Chen et al., 2022). Non-tonal phonemic units (e.g. vowels, consonants) can be defined primarily by **non-pitch articulatory cues**, such as vowel height, voicing, and duration. In contrast, tonal units make use of **pitch cues**, with F0 (fundamental frequency) contour usually considered to be the primary cue (Rhee et al., 2021). In ambiguous contexts, other pitch cues can be used in combination with non-pitch cues such as amplitude, voice quality (e.g. breathy vs creaky), and spectral tilt (Rhee et al., 2021). The Tone Bearing Unit (TBU) varies across languages, with some bearing it on all morphemes, whilst others demonstrate TBU only on specific contrasts or lexical pairs e.g. in Japanese (Jun and Kubozono, 2020).

We compare SLMs trained on non-tonal language as well as three fully lexical tonal languages: Mandarin, Cantonese and Vietnamese. The models are tested primarily on Mandarin data. Mandarin demonstrates full tonality, with tone found on each morpheme (Hyman, 2018), and has been widely studied for tone perception as well as acquisition. Secondly, we also test on data from another lexical tone language, Vietnamese, to assess if our results generalize.

**Mandarin Chinese** is typically described as containing four (lexical) tones and one neutral tone that only occurs in unstressed syllables (Wu et al., 2020). The tones are conventionally assigned the labels 1-4 (T1-4); Figure 1 illustrates the four Mandarin tones. The TBU in Mandarin is morphemic; that is, each morpheme contains one tonal unit. Since one morpheme (one character) corresponds to one tone in Mandarin Chinese, we can use the Pinyin transcription to obtain our tone labels easily (see Figure 1 for notations); for example:

- (1) 今天天气很好  
 Jīn tiān tiān qì hěn hǎo  
 T1 T1 T1 T4 T3 T3

‘The weather today is very good.’

The tone label corresponds to the tone of the character when it is pronounced in isolation (base form). However, Mandarin features *tone sandhi*, i.e. the tone assigned to individual morphemes can change in pronunciation based on the tone of the adjacent

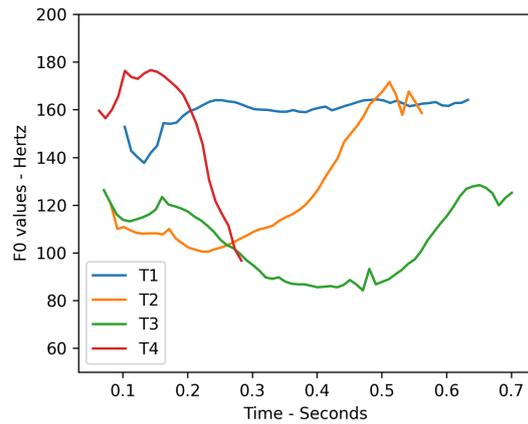


Figure 1: F0 contours of the four Mandarin tones measured from pronunciations recorded by one of the co-authors, a native speaker of Mandarin Chinese. The four syllables are pronounced in isolation (notation: mā T1, má T2, mǎ T3, mà T4).

morpheme (sandhi form). One instance of tone sandhi rules in Mandarin is T3 sandhi (Chen, 2000): if two T3 (‘dipping’) tones occur next to one another, the first will adjust to T2 (‘rising’) to avoid two consecutive T3 tones, as can be seen in examples 2 and 3, after Chen (2000). Tone labels obtained from Pinyin transcriptions only take the base form into account.

- (2) 小  
 xiǎo  
 T3  
 ‘small’

- (3) 小 狗  
 xiǎo gǒu  
 T3 T3 *base form*  
 T2 T3 *sandhi form*  
 ‘small dog’

The primary pitch cue that distinguishes the individual Mandarin tones from each other is F0; however, secondary pitch cues are also present such as voice quality and spectral tilt (Belotel-Grenie and Grenie, 1994; Huang, 2020).

**Vietnamese** also has obligatory tones on every syllable, similar to Mandarin’s morphemic TBU (Kirby, 2011). We adhere to the eight tone system described by Kirby (2011) in our experiment setup.

**Cantonese** is a Sinitic language related to Mandarin, and also features lexical tone, with six tonal distinctions (Zee, 1991) as opposed to Mandarin’s four.

### 3 Related work

The present paper builds both on works interpreting the inner workings of SLMs and on experiments on perception of aspects of human speech.

#### 3.1 Analyzing SLMs

The transformer architecture (Vaswani et al., 2017) has dominated the SLM realm. Researchers have developed many methods to analyze the inner-working of these models. Pasad et al. (2021) provide an overview on the variety of linguistics features encoded by self-supervised SLMs. The models tend to follow an autoencoder-like behavior with the middle layers showing the strongest encoding of a variety of linguistic features.

More recently, research has focused on specific properties of the input audio that is being encoded by the models. Martin et al. (2023) tested whether SLMs can distinguish between voiced and voiceless consonants. Shen et al. (2023) showed that self-supervised as well as visually-supervised SLMs are capable of encoding syntactic properties to some extent. Some prior works in the field have touched on the encoding of suprasegmental features in SSL speech models. Bartelds et al. (2022) showed the hidden state activations of SLMs are capable of capturing intonational and durational information on the phrase level, indicating that they can encode non-segmental information to a significant degree.

Many recent interpretability studies are inspired by psycholinguistics and child language development research. With the rise of probing and other interpretability methods, researchers replicated experimental paradigms in psychology and linguistics to better understand the capabilities of models compared to humans. For example, Wilcox et al. (2023) tested text language models using psycholinguistic experimental paradigms, showing that they are capable of learning syntactic dependencies with relatively little input data.

On the speech side, Lavechin et al. (2023) presented evidence that self-supervised SLMs can develop limited language-specific perception. Cruz Blandón et al. (2023) proposed comparing model behavior using checkpoints in the SLM pre-training process with data in child language development. They showed that computational language models can be a valuable resource in testing or confirming linguistic theories in the language development field. The methodology mostly concerns of the overall learning of the language model

in the output stage. Our work contributes to the explanation of the inner workings of SLMs.

#### 3.2 Human perception experiments

In terms of tone perception, F0 is a clear primary cue (Ryant et al., 2014b; Rhee et al., 2021; Chen et al., 2022), but other secondary pitch cues serve to assist when speech is ambiguous and/or disrupted. Given that conversational speech contains non-trivial speech recognition difficulties such as e.g. tone sandhi and coarticulation, individual variation, and context omission (Ryant et al., 2014b), secondary cues play a role in the distinction of tones. An example of this is voice quality, where for example lowering F0 (introducing ‘creaky’ voice) increased perceptual saliency for T3, whereas T1 and T4 accuracy decreased and T2 remained unaffected (Huang, 2020; Chai, 2019; Kuang, 2017). This emphasises the fact that F0 does not operate in isolation, but that covariation between pitch and voice quality is inherent in Mandarin. Spectral cues (e.g. amplitude differences, spectral tilt) have also been suggested to be sufficient for adult speakers in tone production, while children are thought to hyperarticulate the tonal differences in speech (Rhee et al., 2021).

Suprasegmental cues appear to be preferred in experiments that compare segmental and suprasegmental cues against each other. Human infants are more sensitive to suprasegmental cues, with even newborns showing the same preference (Mehler et al., 1988; Nazzi et al., 1998). Several studies observe that tonal sensitivity develops earlier than perception of vowels and consonants (Xi et al., 2009; Yeung et al., 2013), with sensitivity to non-native tonal distinctions remaining longer than perception of non-tonal non-native phoneme categories (Liu and Kager, 2014; Shi et al., 2017).

Comparing vowels, consonants and tones, Singh et al. (2015) show that Mandarin learning children’s sensitivity to consonants and vowels develop at a similar rate and shows departure from tones. The effect of tone mispronunciation is much larger than that of vowel or consonant mispronunciation for toddlers, but the pattern is reversed in preschoolers (Singh et al., 2015).

#### 3.3 Automatic classification of tones

Automatic tone classification in Mandarin traditionally uses F0 contour and mel-frequency cepstral coefficients (MFCC) features. Advances in deep learning brought improvements in performance of

tone classification. [Ryant et al. \(2014a\)](#) compare MFCC features and F0 contour as input to a neural tone classifier. MFCC features, while not explicitly encoding the F0 contour information, achieve an error rate of 15.56% for Tone 1-4 classification. The combination of MFCC features and F0 contours extracted with different methods did not see an improvement in the classifier’s performance, indicating that the classifier was able to extract F0 contour from the MFCC features, or it was able to predict Mandarin tones reliably without F0 contour information. However, it is possible that the classifier was able to exploit associations between specific phonemes strings and tone labels, and hence avoid learning to detect tone based on pitch and voice quality cues.

After the introduction of self-supervised SLMs, [Yuan et al. \(2021\)](#) fine-tuned an English pre-trained wav2vec2 model ([Baevski et al., 2020](#)) for Mandarin tone classification and achieved a tone error rate of 6% on the same dataset as ([Ryant et al., 2014a](#)). Clearly, SLMs can handle the task of classifying Mandarin lexical tone with labeled fine-tuning. The aim of the present paper is not to compete with the existing implementations of Mandarin tone classifiers; rather we aim to uncover the tone encoding capabilities emerged without explicit supervision.

## 4 Methodology

We use a number of wav2vec2-based ([Baevski et al., 2020](#)) models pre-trained and fine-tuned on datasets of different languages for our investigation. As examples of tonal languages, we choose Mandarin, Vietnamese and Cantonese, whereas English and French serve as non-tonal language examples. The models trained in the languages above are then tested on test data from Mandarin and Vietnamese.

To examine the encoding of tone, we train linear probing classifiers on the hidden state activations extracted from the aforementioned models for every morpheme in our testing datasets.

### 4.1 Datasets

**Training data.** We examine SLMs that were trained on datasets of the following languages:

**Mandarin** pre-trained with AISHELL-2 ([Du et al., 2018](#)) and fine-tuned with AISHELL-1 ([Bu et al., 2017](#)). **English** pre-trained and fine-tuned with LibriSpeech ([Panayotov et al., 2015](#)). **Vietnamese** pre-trained with unlabelled YouTube au-

dio and fine-tuned with the VLSP dataset for ASR ([Nguyen, 2021](#)). **Cantonese** pre-trained on a combined dataset of older Cantonese adult speech and YouTube audio ([Huang and Mak, 2023](#)). **French** pre-trained on MLS French ([Pratap et al., 2020](#)). Table 1 summarizes the characteristics of these datasets.

**Test data.** We primarily use the Mandarin Chinese THCHS-30 dataset ([Wang and Zhang, 2015](#)) for testing models’ encoding of Mandarin tone. THCHS-30 consists of 30 hours of Mandarin speech recorded in a laboratory environment. The dataset is transcribed into both Chinese characters and Mandarin Pinyin. We also obtain character-level forced alignment with the Charsiu aligner ([Zhu et al., 2022](#)).

To test the generalizability of our results, we also use the Vietnamese VIVOS dataset ([Luong and Vu, 2016](#)), which consists of 15 hours of Vietnamese read speech recorded in a laboratory environment. The dataset is transcribed into Vietnamese orthography. We then convert the transcription into International Phonetic Alphabet (IPA) with tone labels with vPhon ([Kirby, 2008](#)). We use the Montreal Forced Aligner ([McAuliffe et al., 2017](#)) to obtain a syllable-level forced alignment.

**Pre-training data.** For the experiments on SLM’s learning trajectory and perceptual patterns (see Section 5.3), we pre-train SLMs from scratch on the following datasets:

- MAGICDATA ([Magic Data Technology Co., 2019](#)), containing 755 hours of read Mandarin Chinese. The dataset was pre-split into a 712-hour training set and a 28-hour validation set.
- LibriSpeech ([Panayotov et al., 2015](#)), see details in Table 1. We split a subset of the LibriSpeech dataset into a 710-hour training set and a 29-hour validation set.

### 4.2 Spoken Language Models

**Architecture.** With the exception of the Cantonese model, all models investigated in this paper are based on the base configuration of wav2vec2 ([Baevski et al., 2020](#)). Wav2vec2-base consists of five convolutional feature encoder and twelve transformer layers. The feature encoder processes the audio waveform input into latent speech representations, and the transformer layers encode the feature encoder output into contextual representations. The wav2vec2-base models has 95M parameters. The

Training language	Tonality	Size (hours)		Speech type
		Pre-training	Fine-tuning	
English (Baevski et al., 2020)	Non-tonal	960	960	Read
French (Parcollet et al., 2023)	Non-tonal	1,000	-	Read
Mandarin (Lu and Chen, 2022)	Tonal	1,000	178	Read
Vietnamese (Nguyen, 2021)	Tonal	13,000	250	YouTube audio/Read
Cantonese (Huang and Mak, 2023)	Tonal	2,800	-	Spontaneous + Read

Table 1: Description of the datasets used in pre-training/fine-tuning models.

Cantonese model uses the wav2vec2-conformer architecture with 180M parameters.

**Training objectives.** The fully self-supervised *pre-training* objective in wav2vec2 consists in discriminating between the matched and unmatched segment representations for a masked portion of the latent speech representation. The ASR *fine-tuning* objective consists in transcribing the audio input into output tokens in the orthography of the target language and is realized by adding a linear layer on top of a pre-trained wav2vec2 model.

**Checkpoints.** For the experiments in Section 5.3 we pre-train two SLMs with the fairseq toolkit (Ott et al., 2019) on LibriSpeech for English and MAG-ICDATA for Mandarin; we train both models for 85,000 steps using 8 Nvidia A100-40GB GPU with update frequency = 8 to simulate training with 64 GPUs. Each model finished training in approximately 96 hours. We save checkpoints every 5,000 steps.

### 4.3 Probing classifiers

**Preprocessing.** We follow previous work (Ryant et al., 2014a) in removing segments transcribed with the neutral tone from the Mandarin tone classification task. Mandarin neutral tones primarily appear in unstressed syllables (cf. Section 2) and hence are more susceptible to variations.

**Generating classifier input.** We extract the hidden state activations of models as a response to audio samples in the test data. We average-pool the hidden state output corresponding to the duration of individual syllables to obtain a vector using forced alignment timestamps. The resulting 768-dimensional vectors are input to the classifiers. To control for the influence of lexical cues on tone detection, we construct an exclusive train-test-split such that phoneme strings appearing in the test set do not appear in the training set. This setup

Language	Split	Samples
Mandarin	Train	223,851
Mandarin	Test	45,772
Vietnamese	Train	124,248
Vietnamese	Test	29,629

Table 2: Train/test splits for the tone probing classifier, for the Mandarin and Vietnamese data.

prevents the probing classifier from exploiting associations between tones and phoneme sequences. We employ a randomized 80:20 train-test split with the split sizes shown in Table 2.

**F0 and MFCC baselines.** We closely follow Ryant et al. (2014a) and use F0 contours and 40-dimensional mel-frequency cepstral coefficients (MFCC) features as baselines. We use Librosa (McFee et al., 2023) to extract the MFCC features and Praat (Jadoul et al., 2018; Boersma and Weenink, 2021) to extract the F0 contours from the audio samples. We then find the center frame for each word using the alignment timestamps and concatenate all frames in a 21 frame window (10-1-10) for both F0 and MFCC features. We end up with a 21-dimensional vector for F0 contours and 840-dimensional vector for MFCC features as our baseline classifier inputs.

**Text baseline.** In addition to audio baselines, we also include a text-based transformer model in our comparison. BERT (Devlin et al., 2019) serves as a reference point to show how much information is encoded in the speech signal as opposed to what can be guessed from pure text. We use a Chinese pre-trained BERT<sup>1</sup> that encodes Chinese characters into vectors. We extract per-word hidden state outputs with a resulting 768-dimensional vector.

<sup>1</sup><https://huggingface.co/bert-base-chinese>

Split	Samples
Train	92,413
Test	15,688

Table 3: Train/test split for the consonant probing classifier, for the Mandarin data.

**Tone classifiers.** We use the syllable activation vectors as input to a Ridge linear classifier that predicts the lexical tone of the input morpheme. We select the final model via 5-fold cross-validation, and report the classification accuracy on the test split. The regularization strength  $\alpha$  was tuned for values  $\{10^n \mid n \in \{-4, -3, -2, -1, 0, 1, 2\}\}$ .

**Consonant classifiers.** When comparing tone to consonant classification, we employ the same classifier setup for consonant and replicate the perception experiment in Wang and Chen (2020) in Section 5.3.1. Since we only investigate consonants that appear solely in the onset position and the rest of the phonemes are not relevant to our task, we use the same syllable vectors as above instead of obtaining a phoneme vector with using phoneme level alignment. We construct exclusive train-test splits that contain unique rhymes (nucleus + coda) of the syllables. Specific details of the train/test split for this experiment can be found in Table 3.

## 5 Results

In this section, we present a series of experiments for analyzing the encoding of tone in SLMs.

### 5.1 Tone encoding across languages

Figure 2 shows the tone classification accuracy using the layer-wise representations of all models pre-trained on non-tonal (left) versus tonal (panel) languages. We see that all layers of all models perform better than the F0 and MFCC baselines, which themselves outperform the text-based BERT baseline. The classification accuracy for tonal language models is overall higher, and increases in the higher layers of the models. Models trained on non-tonal languages also show substantial encoding of tone; but remarkably, there is a substantial drop in classification accuracy in their final layers while the corresponding decrease is much less pronounced in tonal language models.

We repeat the tone classification experiment for Vietnamese tones. Results in Figure 3 show the Cantonese model performs slightly better than the

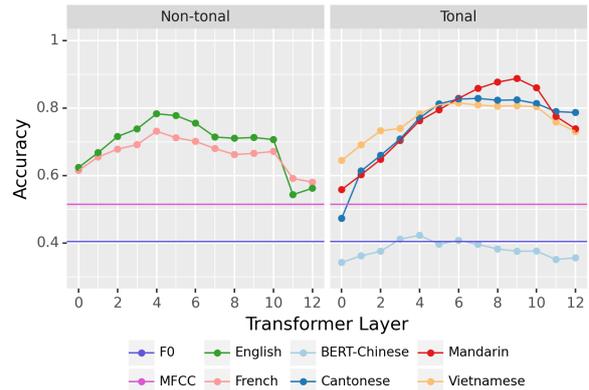


Figure 2: Classification accuracy of Mandarin lexical tones using layer-wise representations from models pre-trained on tonal and non-tonal languages.

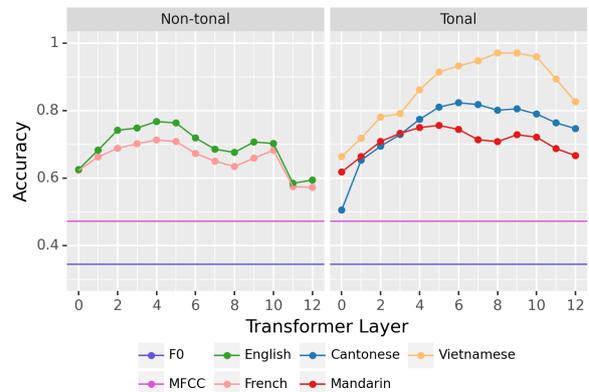


Figure 3: Classification accuracy of Vietnamese lexical tones with hidden-state activations from models pre-trained on tonal and non-tonal languages.

English model, especially towards the later layers; the Mandarin model, however, patterns similar to the English model. This is likely due to the fact that Mandarin has fewer tonal contrasts than Vietnamese and Cantonese (cf. Section 2).

Studies on human participants show that speakers of other tonal languages perform better at identifying Mandarin lexical tones compared to non-tonal language speakers (So and Best, 2010), and the SLMs we tested show the same pattern. Regarding Vietnamese tones, the result is more equivocal suggesting that Cantonese tone representations generalize to Vietnamese to some extent, while Mandarin ones do not.

### 5.2 Impact of ASR fine-tuning

We examine how fine-tuning for ASR impacts the encoding of tone in SLMs. Since tonal information is crucial for correctly transcribing tonal language input, SLMs fine-tuned for tonal languages are ex-

476 expected to perform better at our tone classification  
 477 task. Figure 4 compares the tone classification  
 478 accuracy of the English and Mandarin models, pre-  
 479 trained only (left) versus pre-trained and then fine-  
 480 tuned (right); Figure 5 shows the corresponding  
 481 results for English and Vietnamese.

482 We find that fine-tuning affects the encoding  
 483 of tone for non-tonal vs tonal language models  
 484 in opposite ways: classification accuracy benefits  
 485 from fine-tuning for Mandarin, but is harmed by  
 486 it for the English model. The same pattern holds  
 487 for English vs Vietnamese on the Vietnamese tone  
 488 data in Figure 5.

489 These results likely reflect the fact that ASR  
 490 fine-tuning encourages the SLM to increase its spe-  
 491 cialization in identifying the language-specific in-  
 492 formation needed to output the written form of the  
 493 language. Tonal information may not contribute  
 494 much to this objective in non-tonal languages, and  
 495 thus fine-tuning would tend to remove it. In tonal-  
 496 language ASR however, tone information may be  
 497 crucial to correctly transcribe the input audio, for  
 498 example, when disambiguating Mandarin syllables  
 499 that consist of the same phonemes and only dif-  
 500 fer in tone, in order to output the correct Chinese  
 501 character.

### 5.3 Comparison to human perception

502 In this section we report the results motivated by  
 503 tone and consonant perception patterns in humans.  
 504

#### 5.3.1 Learning trajectory

505 Children have a higher sensitivity to tone than con-  
 506 sonant distinctions early on. For children speak-  
 507 ing a non-tonal language, this sensitivity towards  
 508 tone continues longer than sensitivity towards non-  
 509 native segmental features, i.e. consonants and vowels  
 510 (Shi et al., 2017; Liu and Kager, 2014).  
 511

512 Here we aim to determine the corresponding  
 513 learning trajectory in SLMs by testing them during  
 514 pre-training. Figure 6 shows the accuracy of clas-  
 515 sifying Mandarin consonants and tones in the best  
 516 performing layer of SLMs trained on English and  
 517 Mandarin as a function of the number of training  
 518 steps.

519 Although we observe classification accuracy of  
 520 the SLMs quickly surpasses the F0 and MFCC  
 521 baselines after 10,000 steps, we do not detect an  
 522 obvious difference in the overall pattern between  
 523 the case of consonants and tones. This suggest  
 524 that SLMs do not follow the same differential tra-  
 525 jectory as children, at least as measured via our

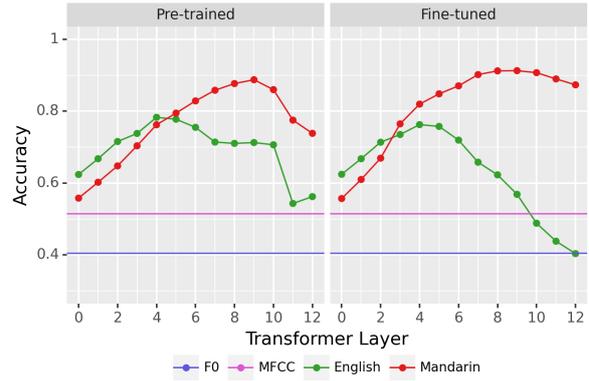


Figure 4: Classification accuracy of Mandarin lexical tones using layer-wise representations from models pre-trained and fine-tuned on Mandarin and English.



Figure 5: Classification accuracy of Vietnamese lexical tones using layer-wise representations from models pre-trained and fine-tuned on Vietnamese and English.

methodology.

526

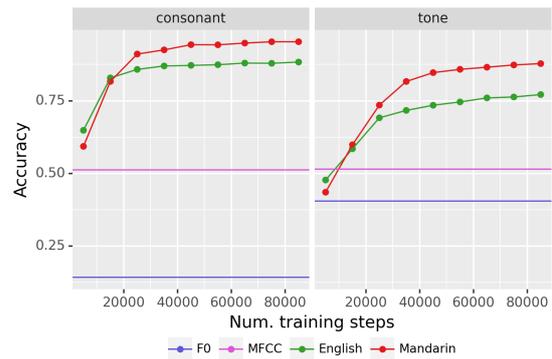


Figure 6: Classification accuracy of Mandarin lexical tones versus consonants for models pre-trained on English and Mandarin.

#### 5.3.2 Tone and consonant contrasts

527

528 Non-native speakers can have difficulty distinguish-  
 529 ing between T2-T3 and T1-T4 tone pairs in Man-  
 530 darin (Hao, 2012). We investigate this pattern in

528

529

530

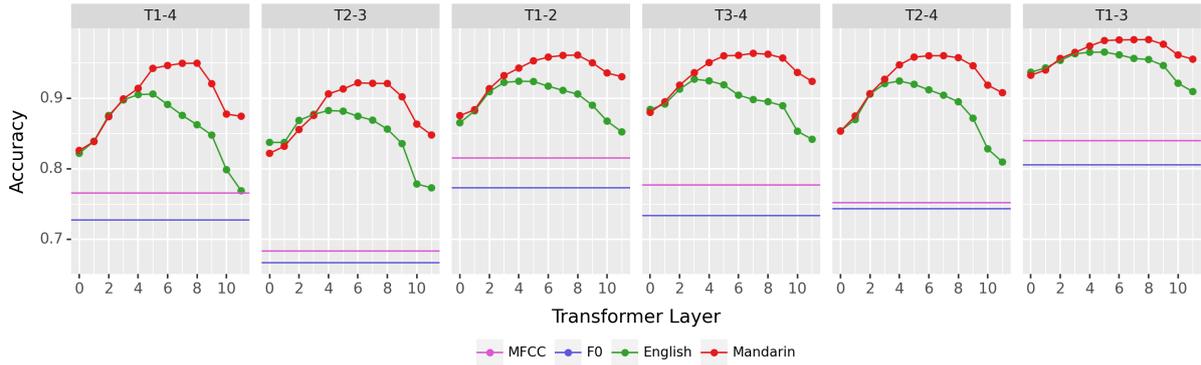


Figure 7: Binary classification accuracy for Mandarin tonal pairs, for English and Mandarin models.

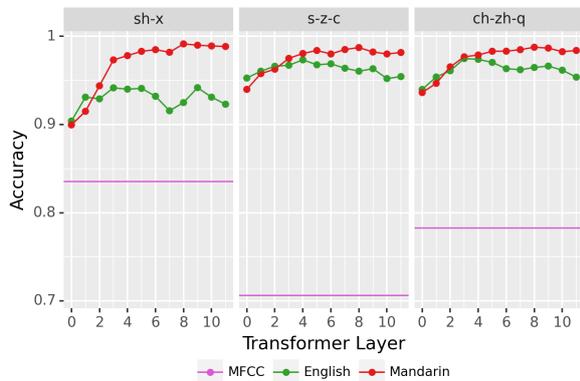


Figure 8: Classification accuracy for Mandarin consonant groups, for English and Mandarin models. The F0 baseline with its much lower classification accuracy is omitted from this figure for clarity.

Group	Mandarin		English	
	Pinyin	IPA	Alphabet	IPA
1	sh, x	ʃ, ç	sh	ʃ
2	ch, zh, q	tʃ <sup>h</sup> , tʂ, tʃ <sup>h</sup>	ch	tʃ
3	s, z, c	s, ts, ts <sup>h</sup>	s	s

Table 4: Perceptual mapping of Mandarin consonants onto English consonants (Wang and Chen, 2020).

by Wang and Chen (2020).

## 6 Conclusion

We analyze the tone encoding capabilities of spoken language models trained on three tonal and two non-tonal languages, using classifier probes with data from two tonal languages Mandarin and Vietnamese. We find that SLMs trained on either tonal or non-tonal languages encode tonal information to a significant degree.

We also find that fine-tuning for the speech recognition task enhances the tone encoding capabilities of models trained on tonal languages but reduces them for models trained on non-tonal languages. While we see evidence suggesting that the learning trajectories of SLMs in pre-training do not follow the same developmental trajectories found in human language acquisition, we find that SLMs show patterns similar to that of human listeners in tone and consonant perception experiments.

Here we investigate the encoding of tone, one example of suprasegmentals; encoding of other suprasegmental features such as stress patterns and intonation is important to study in future work.

pre-trained SLMs via a dedicated probing experiment, using the final (85,000 steps) checkpoint of the pre-trained models in Section 5.3.1. As can be seen in Figure 7, tone pairs T1-4 and T2-3 show the largest differences in the best classification accuracy between the Mandarin and English models, which roughly matches human perceptual pattern.

We complement the results on the development of tone contrast with a parallel experiment on those Mandarin consonant contrasts which are challenging for speakers of English. Each member of a contrasting group is perceived as the same phoneme by English speakers due to perceptual assimilation (Wang and Chen, 2020). Table 4 displays the resulting mapping to English phoneme categories.

Figure 8 shows that accuracy for consonant groups 2 and 3 match closely for the two models. Group 1 shows a discrepancy, possibly due to the potential mapping of Mandarin x /ç/ into two English consonants sh /ʃ/ and z /z/, as hypothesized

## 7 Limitations

We selected SLMs based on the wav2vec2 architecture in our experimental design, but we acknowledge that the training data of the models selected is quite varied in their size and quality (noisy vs clean speech) as described in Section 4.1. This is partially due to the scarce availability of (high-quality) speech data for underrepresented languages, especially the many tonal languages of the world. Hence SLMs pre-trained on monolingual datasets of these languages are also sparse. The Vietnamese wav2vec2 model (Nguyen, 2021) was trained on a significantly larger amount of data (13k hours) than the other models tested (around 1000 hours). It is possible that in addition to the inclusion of tonal languages in training, the amount of training data also played a role in increasing the tonal encoding capabilities of SLMs. However, literature has shown that more training data does not always have a positive impact on the models performance if the additional data is noisy (Parcollet et al., 2023). At the same time, we note that the Cantonese model (Huang and Mak, 2023), in addition to being pre-trained on a larger dataset, is also different in architecture. Additionally, our use of two read speech datasets as test data does not fully reflect the linguistic diversity of different accents and dialects in Mandarin and Vietnamese. Future work needs to go wider and deeper in both model architecture and dataset diversity in order to uncover more generalizable patterns in different languages.

## References

Badr M. Abdullah, Iuliia Zaitova, Tania Avgustinova, Bernd Möbius, and Dietrich Klakow. 2021. How Familiar Does That Sound? Cross-Lingual Representational Similarity Analysis of Acoustic Word Embeddings.

Alexei Baevski, Henry Zhou, Abdelrahman Mohamed, and Michael Auli. 2020. Wav2vec 2.0: A Framework for Self-Supervised Learning of Speech Representations.

Martijn Bartelds, Wietse de Vries, Faraz Sanal, Caitlin Richter, Mark Liberman, and Martijn Wieling. 2022. Neural representations for modeling variation in speech. *Journal of Phonetics*, 92:101137.

Agnes Belotel-Grenie and Michel Grenie. 1994. Phonation types analysis in standard chinese. In *3rd International Conference on Spoken Language Processing, ICSLP 1994*, pages 343–346. The Interna-

tional Society for Computers and Their Applications (ISCA).

Paul Boersma and David Weenink. 2021. Praat: Doing phonetics by computer [Computer program].

Hui Bu, Jiayu Du, Xingyu Na, Bengu Wu, and Hao Zheng. 2017. AISHELL-1: An Open-Source Mandarin Speech Corpus and A Speech Recognition Baseline.

Yuan Chai. 2019. THE SOURCE OF CREAK IN MANDARIN UTTERANCES.

Matthew Y. Chen. 2000. *Tone Sandhi: Patterns across Chinese Dialects*. Cambridge Studies in Linguistics. Cambridge University Press, Cambridge.

Yue Chen, Yingming Gao, and Yi Xu. 2022. Computational Modelling of Tone Perception Based on Direct Processing of f0 Contours. *Brain Sciences*, 12(3):337.

Alexis Conneau, German Kruszewski, Guillaume Lample, Loïc Barrault, and Marco Baroni. 2018. What you can cram into a single &#!#\* vector: Probing sentence embeddings for linguistic properties. In *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 2126–2136, Melbourne, Australia. Association for Computational Linguistics.

María Andrea Cruz Blandón, Alejandrina Cristia, and Okko Räsänen. 2023. Introducing Meta-analysis in the Evaluation of Computational Models of Infant Language Development. *Cognitive Science*, 47(7):e13307.

Maureen de Seyssel, Marvin Lavechin, Yossi Adi, Emmanuel Dupoux, and Guillaume Wisniewski. 2022. Probing phoneme, language and speaker information in unsupervised speech representations. In *Interspeech*.

Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding.

Jiayu Du, Xingyu Na, Xuechen Liu, and Hui Bu. 2018. AISHELL-2: Transforming Mandarin ASR Research Into Industrial Scale.

Yen-Chen Hao. 2012. Second language acquisition of Mandarin Chinese tones by tonal and non-tonal language speakers. *Journal of Phonetics*, 40(2):269–279.

John Hewitt and Christopher D. Manning. 2019. A Structural Probe for Finding Syntax in Word Representations. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pages 4129–4138, Minneapolis, Minnesota. Association for Computational Linguistics.

678	Wei-Ning Hsu, Benjamin Bolte, Yao-Hung Hubert	Danni Ma, Neville Ryant, and Mark Liberman. 2021.	730
679	Tsai, Kushal Lakhota, Ruslan Salakhutdinov, and	<a href="#">Probing Acoustic Representations for Phonetic Prop-</a>	731
680	Abdelrahman Mohamed. 2021. <a href="#">HuBERT: Self-</a>	<a href="#">erties</a> . In <i>ICASSP 2021 - 2021 IEEE International</i>	732
681	<a href="#">Supervised Speech Representation Learning by</a>	<i>Conference on Acoustics, Speech and Signal Pro-</i>	733
682	<a href="#">Masked Prediction of Hidden Units</a> .	<i>cessing (ICASSP)</i> , pages 311–315.	734
683	Ranzo Huang and Brian Mak. 2023. <a href="#">Wav2vec 2.0 ASR</a>	Ltd. Magic Data Technology Co. 2019. <a href="#">MAGICDATA</a>	735
684	<a href="#">for Cantonese-Speaking Older Adults in a Clinical</a>	<a href="#">Mandarin Chinese Read Speech Corpus</a> .	736
685	<a href="#">Setting</a> . In <i>INTERSPEECH 2023</i> , pages 4958–4962.	Kinan Martin, Jon Gauthier, Canaan Breiss, and Roger	737
686	ISCA.	Levy. 2023. <a href="#">Probing Self-supervised Speech Mod-</a>	738
687	Yaqian Huang. 2020. <a href="#">Different attributes of creaky</a>	<a href="#">els for Phonetic and Phonemic Information: A Case</a>	739
688	<a href="#">voice distinctly affect Mandarin tonal perception</a> .	<a href="#">Study in Aspiration</a> . In <i>INTERSPEECH 2023</i> , pages	740
689	<i>The Journal of the Acoustical Society of America</i> ,	251–255. ISCA.	741
690	147(3):1441–1458.	Michael McAuliffe, Michaela Socolof, Sarah Mihuc,	742
691	Larry M. Hyman. 2018. What tone teaches us about	Michael Wagner, and Morgan Sonderegger. 2017.	743
692	language. <i>Language</i> , 94(3):698–709.	<a href="#">Montreal Forced Aligner: Trainable Text-Speech</a>	744
693	Yannick Jadoul, Bill Thompson, and Bart de Boer.	<a href="#">Alignment Using Kaldi</a> . In <i>Interspeech 2017</i> , pages	745
694	2018. <a href="#">Introducing Parselmouth: A Python interface</a>	498–502. ISCA.	746
695	<a href="#">to Praat</a> . <i>Journal of Phonetics</i> , 71:1–15.	Brian McFee, Matt McVicar, Daniel Faronbi, Iran Ro-	747
696	Sun-Ah Jun and Haruo Kubozono. 2020. Asian Pacific	man, Matan Gover, Stefan Balke, Scott Seyfarth, Ay-	748
697	Rim. In Carlos Gussenhoven and Aoju Chen, ed-	oub Malek, Colin Raffel, Vincent Lostanlen, Ben-	749
698	itors, <i>The Oxford Handbook of Language Prosody</i> ,	jamin Van Niekirk, Dana Lee, Frank Cwitkowitz,	750
699	Oxford Handbooks. Oxford University Press, Ox-	Frank Zalkow, Oriol Nieto, Dan Ellis, Jack Ma-	751
700	ford, New York.	son, Kyungyun Lee, Bea Steers, Emily Halvachs,	752
701	James Kirby. 2008. vPhon: A Vietnamese phonetizer	Carl Thomé, Fabian Robert-Stöter, Rachel Bittner,	753
702	(version 2.1.1).	Ziyao Wei, Adam Weiss, Eric Battenberg, Keun-	754
703	James P. Kirby. 2011. <a href="#">Vietnamese (hanoi vietnamese)</a> .	woo Choi, Ryuichi Yamamoto, CJ Carr, Alex Met-	755
704	<i>Journal of the international phonetic association</i> ,	sai, Stefan Sullivan, Pius Friesch, Asmitha Krish-	756
705	41(3):381–392.	nakumar, Shunsuke Hidaka, Steve Kowalik, Fabian	757
706	Jianjing Kuang. 2017. <a href="#">Covariation between voice qual-</a>	Keller, Dan Mazur, Alexandre Chabot-Leclerc,	758
707	<a href="#">ity and pitch: Revisiting the case of Mandarin creaky</a>	Curtis Hawthorne, Chandrashekhara Ramaprasad,	759
708	<a href="#">voice</a> . <i>The Journal of the Acoustical Society of</i>	Myungchul Keum, Juanita Gomez, Will Mon-	760
709	<i>America</i> , 142(3):1693–1706.	roe, Viktor Andreevitch Morozov, Kian Eliasi,	761
710	Marvin Lavechin, Maureen De Seyssel, Marianne Mé-	Nullmightybofo, Paul Biberstein, N. Dorukhan Ser-	762
711	tais, Florian Metze, Abdelrahman Mohamed, Hervé	gin, Romain Hennequin, Rimvydas Naktinis, Bean-	763
712	Bredin, Emmanuel Dupoux, and Alejandrina Cristia.	towel, Taewoon Kim, Jon Petter Åsen, Joon Lim,	764
713	2023. Statistical learning models of early phonetic	Alex Malins, Darío Hereñú, Stef Van Der Stru-	765
714	acquisition struggle with child-centered audio data.	ijk, Lorenz Nickel, Jackie Wu, Zhen Wang, Tim	766
715	Liquan Liu and René Kager. 2014. <a href="#">Perception of tones</a>	Gates, Matt Vollrath, Andy Sarroff, Xiao-Ming,	767
716	<a href="#">by infants learning a non-tone language</a> . <i>Cognition</i> ,	Alastair Porter, Seth Kranzler, VoodooHop, Mat-	768
717	133(2):385–394.	tia Di Gangi, Helmi Jinoz, Connor Guerrero, Ab-	769
718	Ke-Han Lu and Kuan-Yu Chen. 2022. <a href="#">A context-</a>	duttayyeb Mazhar, Toddme2178, Zvi Baratz, An-	770
719	<a href="#">aware knowledge transferring strategy for CTC-</a>	ton Kostin, Xinlu Zhuang, Cash TingHin Lo, Pavel	771
720	<a href="#">based ASR</a> .	Campr, Eric Semeniuc, Monsij Biswal, Shayenne	772
721	Hieu-Thi Luong and Hai-Quan Vu. 2016. A non-expert	Moura, Paul Brossier, Hojin Lee, and Waldir Pi-	773
722	Kaldi recipe for Vietnamese speech recognition sys-	menta. 2023. <a href="#">Librosa/librosa: 0.10.1</a> . Zenodo.	774
723	tem. In <i>Proceedings of the Third International Work-</i>	J. Mehler, P. Jusczyk, G. Lambertz, N. Halsted,	775
724	<i>shop on Worldwide Language Service Infrastructure</i>	J. Bertoncini, and C. Amiel-Tison. 1988. <a href="#">A precu-</a>	776
725	<i>and Second Workshop on Open Infrastructures and</i>	<a href="#">ursor of language acquisition in young infants</a> . <i>Cogni-</i>	777
726	<i>Analysis Frameworks for Human Language Tech-</i>	<i>tion</i> , 29(2):143–178.	778
727	<i>nologies (WLSI/OIAF4HLT2016)</i> , pages 51–55, Os-	T. Nazzi, J. Bertoncini, and J. Mehler. 1998. <a href="#">Language</a>	779
728	saka, Japan. The COLING 2016 Organizing Commit-	<a href="#">discrimination by newborns: Toward an understand-</a>	780
729	tee.	<a href="#">ing of the role of rhythm</a> . <i>Journal of Experimen-</i>	781
		<i>tial Psychology. Human Perception and Performance</i> ,	782
		24(3):756–766.	783
		Thai Binh Nguyen. 2021. <a href="#">Vietnamese end-to-end</a>	784
		<a href="#">speech recognition using wav2vec 2.0</a> .	785
		Myle Ott, Sergey Edunov, Alexei Baevski, Angela	786
		Fan, Sam Gross, Nathan Ng, David Grangier, and	787

788	Michael Auli. 2019. <a href="#">Fairseq: A Fast, Extensible Toolkit for Sequence Modeling</a> . pages 48–53, Minneapolis, Minnesota. Association for Computational Linguistics.	842
789		843
790		844
791		
792	Vassil Panayotov, Guoguo Chen, Daniel Povey, and Sanjeev Khudanpur. 2015. <a href="#">Librispeech: An ASR corpus based on public domain audio books</a> . In <i>2015 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)</i> , pages 5206–5210.	845
793		846
794		847
795		848
796		
797		
798	Titouan Parcollet, Ha Nguyen, Solene Evain, Marceley Zanon Boito, Adrien Pupier, Salima Mdhaffar, Hang Le, Sina Alisamir, Natalia Tomashenko, Marco Dinarelli, Shucong Zhang, Alexandre Allauzen, Maximin Coavoux, Yannick Esteve, Mickael Rouvier, Jerome Goulian, Benjamin Lecouteux, Francois Portet, Solange Rossato, Fabien Ringeval, Didier Schwab, and Laurent Besacier. 2023. <a href="#">LeBenchmark 2.0: A Standardized, Replicable and Enhanced Framework for Self-supervised Representations of French Speech</a> .	849
799		850
800		851
801		852
802		
803		
804		
805		
806		
807		
808		
809	Ankita Pasad, Ju-Chieh Chou, and Karen Livescu. 2021. <a href="#">Layer-wise Analysis of a Self-supervised Speech Representation Model</a> .	853
810		854
811		
812	Vineel Pratap, Qiantong Xu, Anuroop Sriram, Gabriel Synnaeve, and Ronan Collobert. 2020. <a href="#">MLS: A Large-Scale Multilingual Dataset for Speech Research</a> . In <i>Interspeech 2020</i> , pages 2757–2761.	855
813		856
814		857
815		858
816	Nari Rhee, Aoju Chen, and Jianjing Kuang. 2021. <a href="#">Going beyond F0: The acquisition of Mandarin tones</a> . <i>Journal of Child Language</i> , 48(2):387–398.	859
817		860
818		861
819	Neville Ryant, Malcolm Slaney, Mark Liberman, Elizabeth Shriberg, and Jiahong Yuan. 2014a. <a href="#">Highly Accurate Mandarin Tone Classification In The Absence of Pitch Information</a> . In <i>Speech Prosody 2014</i> , pages 673–677. ISCA.	862
820		863
821		864
822		865
823		
824	Neville Ryant, Jiahong Yuan, and Mark Liberman. 2014b. <a href="#">Mandarin tone classification without pitch tracking</a> . In <i>2014 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)</i> , pages 4868–4872, Florence, Italy. IEEE.	866
825		867
826		868
827		869
828		
829	Gaofei Shen, Afra Alishahi, Arianna Bisazza, and Grzegorz Chrupala. 2023. <a href="#">Wave to Syntax: Probing spoken language models for syntax</a> . In <i>INTER-SPEECH 2023</i> , pages 1259–1263.	870
830		871
831		872
832		873
833		874
834	Rushen Shi, Jun Gao, André Achim, and Aijun Li. 2017. <a href="#">Perception and Representation of Lexical Tones in Native Mandarin-Learning Infants and Toddlers</a> . <i>Frontiers in Psychology</i> , 8.	875
835		876
836		
837	Leher Singh and Charlene S. L. Fu. 2016. <a href="#">A New View of Language Development: The Acquisition of Lexical Tone</a> . <i>Child Development</i> , 87(3):834–854.	877
838		878
839		879
840	Leher Singh, Hwee Hwee Goh, and Thilanga D. We-walaarachchi. 2015. <a href="#">Spoken word recognition in early childhood: Comparative effects of vowel, consonant and lexical tone variation</a> . <i>Cognition</i> , 142:1–11.	880
841		881
		882
	Connie K. So and Catherine T. Best. 2010. <a href="#">Cross-language Perception of Non-native Tonal Contrasts: Effects of Native Phonological and Phonetic Influences</a> . <i>Language and speech</i> , 53(Pt 2):273–293.	883
		884
		885
		886
		887
	Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Lukasz Kaiser, and Illia Polosukhin. 2017. <a href="#">Attention Is All You Need</a> .	
	Dong Wang and Xuewei Zhang. 2015. <a href="#">THCHS-30 : A Free Chinese Speech Corpus</a> .	
	Xinchun Wang and Jidong Chen. 2020. <a href="#">The Acquisition of Mandarin Consonants by English Learners: The Relationship between Perception and Production</a> . <i>Languages</i> , 5(2):20.	
	Ethan Gotlieb Wilcox, Richard Futrell, and Roger Levy. 2023. <a href="#">Using Computational Models to Test Syntactic Learnability</a> . <i>Linguistic Inquiry</i> , pages 1–44.	
	Yaru Wu, Martine Adda-Decker, and Lori Lamel. 2020. <a href="#">Mandarin lexical tones: A corpus-based study of word length, syllable position and prosodic position on duration</a> . pages 1908–1912.	
	Jie Xi, Wei JIANG, Linjun Zhang, and Hua Shu. 2009. <a href="#">Categorical Perception of VOT and Lexical Tones in Chinese and the Developmental Course</a> . <i>Acta Psychologica Sinica</i> , 41:572–579.	
	H. Henny Yeung, Ke Heng Chen, and Janet F. Werker. 2013. <a href="#">When does native language input affect phonetic perception? The precocious case of lexical tone</a> . <i>Journal of memory and language</i> , 68(2):123–139.	
	Moira Yip. 2002. <i>Tone</i> . Cambridge Textbooks in Linguistics. Cambridge University Press, Cambridge.	
	Jiahong Yuan, Neville Ryant, Xingyu Cai, Kenneth Church, and Mark Liberman. 2021. <a href="#">Automatic recognition of suprasegmentals in speech</a> .	
	Eric Zee. 1991. <a href="#">Chinese (Hong Kong Cantonese)</a> . <i>Journal of the International Phonetic Association</i> , 21(1):46–48.	
	Jian Zhu, Cong Zhang, and David Jurgens. 2022. <a href="#">Phone-to-audio alignment without text: A Semi-supervised Approach</a> . <i>IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)</i> .	