

Word-specific tonal realizations in Mandarin

Yu-Ying Chuang¹, Melanie J. Bell², Yu-Hsiang Tseng³, and R. Harald Baayen³

¹National Taiwan Normal University

²Anglia Ruskin University, UK

³Eberhard Karl University of Tübingen, Germany

May, 2024

Abstract

The pitch contours of Mandarin two-character words are generally understood as being shaped by the underlying tones of the constituent single-character words, in interaction with articulatory constraints imposed by factors such as speech rate, co-articulation with adjacent tones, segmental make-up, and predictability. This study shows that tonal realization is also partially determined by words' meanings. We first show, on the basis of a Taiwan corpus of spontaneous conversations, using the generalized additive regression model, and focusing on the rise-fall tone pattern, that after controlling for effects of speaker and context, word type is a stronger predictor of pitch realization than all the previously established word-form related predictors combined. Importantly, the addition of information about meaning in context improves prediction accuracy even further. We then proceed to show, using computational modeling with context-specific word embeddings, that token-specific pitch contours predict word type with 50% accuracy on held-out data, and that context-sensitive, token-specific embeddings can predict the shape of pitch contours with 30% accuracy. These accuracies, which are an order of magnitude above chance level, suggest that the relation between words' pitch contours and their meanings are sufficiently strong to be functional for language users. The theoretical implications of these empirical findings are discussed.

keywords: Tone, Mandarin, embeddings, GAMs, form-meaning isomorphy, rise-fall tone pattern, two-syllable words

Acknowledgements This research was funded by the European Research Council, grant SUBLIMINAL (#101054902) awarded to Harald Baayen. The authors are indebted to Dr. Matteo Fasiolo, University of Bristol, for his statistical advice on the application of GAMs to pitch contours.

Author note Correspondence concerning this article should be addressed to Dr. Yu-Ying Chuang, National Taiwan Normal University, e-mail: yuying.chuang@ntnu.edu.tw.

1 Introduction

The tone system of Mandarin Chinese is commonly described as having four lexical tones, namely high level (T1), rising (T2), dipping/low (T3), and falling (T4), plus one neutral or floating tone whose shape depends on the preceding tone (Chao, 1968). While these tonal representations are well established in the relevant literature, and are taught both in Chinese schools and in second language learning classrooms, it is also known that their realizations, i.e. the actual pitch contours of spoken words, can differ significantly from the canonical descriptions.

Xu (2001) distinguishes two sources of tonal variation: involuntary and voluntary. Involuntary variation arises from articulatory constraints posited to be beyond the speaker’s control. Voluntary variation arises from the speaker’s communicative intentions, such as eliciting a question or creating emphasis, where changes in rhythm, accent placement, and intonation can all give rise to substantial modification of tonal realization; these voluntary constraints are intimately related to the syntactic and pragmatic functions of an utterance (Gårding, 1987; Liu and Xu, 2005; Shen, 1989, 1990a; Xu, 1999). In addition, the realization of tone is modulated by sociolinguistic and paralinguistic factors such as gender, dialect, and emotion (Fon and Chiang, 1999; Zhang et al., 2006).

Involuntary articulatory constraints on tone production arise both from the segmental makeup of syllables and from the processes of connected speech. At the syllable level, vowels, onsets, and syllable structure all contribute to tonal variation (Fon and Hsu, 2007; Ho, 1976; Howie, 1974; Whalen and Levitt, 1995; Xu and Xu, 2003). At the utterance level, where a sequence of tones is produced, the realization of a given tone is greatly influenced by its preceding and following tones, leading to tonal co-articulation (Shen, 1990b; Shih, 1988; Xu, 1997). Co-articulated tones usually deviate from their canonical tonal shapes; in extreme cases the original shapes are no longer preserved (Shen, 1989; Shih and Kochanski, 2000) and may be unrecognisable to native speakers (Xu, 1994). Physically speaking, it requires a certain amount of time to raise or lower pitch (Xu and Sun, 2002), and therefore tonal realizations are very dependent on whether speakers have sufficient time to realize a given tonal contour. Under the time pressure of fast speech, tonal targets can usually not be fully realized (Tang and Li, 2020); this leads to significant deviation from the canonical patterns and often results in tonal reduction (Cheng and Xu, 2015). The accepted descriptions of the tones as ‘level’, ‘rising’, ‘dipping/low’, and ‘falling’ are therefore generalizations across considerable variability in the fine detail of their phonetic realizations.

In the extensive literature on Mandarin tones, their semantic function is straightforwardly simple: different tones distinguish between alternative meanings. For instance, 就 *jiù* ‘then’ and 九 *jiǔ* ‘nine’ are differentiated by having a falling and a dipping tone, respectively. However, the very same combinations of segments and tone often realize many other different meanings, as exemplified by 九 *jiǔ* ‘nine’ and 酒 *jiǔ* ‘alcoholic beverage’. The combination of strong phonotactic constraints on syllable structure and a limited number of lexical tones has given rise to widespread homophony and polysemy, often in combination with homography. For instance, 就 *jiù* ‘then’, has a wide range of translation equivalents in English, including *then*, *at once*, *only*, *already*, *to approach*, *to accomplish*, *to suffer*, and *to take advantage of* (<https://www.pleco.com/>, s.v.).

The examples in the preceding paragraph are all monosyllabic words, hence written with a single character; however, in the Chinese Lexical Database (Sun et al., 2018), only 8% of the 48,000 words are monosyllabic. The majority of Mandarin words are disyllabic, written with two characters.¹ The tonal targets of disyllabic words are taken to be subject to the same voluntary and involuntary constraints that govern the realization of tone in monosyllabic words. As a consequence, all disyllabic words sharing, for example, an initial falling tone and a subsequent rising tone are assumed to be realized with the same underlying pitch contour; any differences in how the tones are realized are assumed to be attributable to the involuntary and voluntary processes described above. However, the present study will show that, alongside the known articulatory constraints, there is a previously undocumented close association between the meanings of Mandarin disyllabic words and the realization of their tonal contours.

The basis for our study was laid by a growing body of research on English showing that fine-grained phonetic variation can be systematically associated with differences in meaning. For example, at the word level, Gahl (2008) showed that homophones such as *time* and *thyme* are realized with different acoustic durations. In the same vein, Lohmann (2018) found that the durations of words such as *cut* depend on whether they are used as nouns or verbs. These differences in word duration were initially explained as a consequence of the different relative frequencies of the homophones in these studies. However, Gahl and Baayen (2024) showed that the meanings of English homophones are a strong co-determinant of their spoken word durations even after frequency differences and other co-determinants such as speech rate are taken

¹In this study we use a corpus of Taiwan Mandarin spontaneous speech (Fon, 2004) and take the word labels supplied by the corpus as given. The labeled words include, for example, nouns such as 學校 *xuéxiào* ‘school’, verb forms such as 學到 *xué dào* ‘learn+resultative’, and negated verbs such as 不是 *búshì* ‘not+be’. Since all these forms have two syllables, we refer to them collectively as disyllabic words.

into account. A relationship between meaning and duration has also been found for the English suffix /s/: different grammatical functions of this suffix (e.g., plural and third person singular) tend to be realized with different durations (Plag et al., 2017). Furthermore, the relationship between meaning and phonetic realization may extend beyond durational differences. Drager (2011) reported that the phonetic realization of the word *like* varies according to its discourse or grammatical meanings, not only in the duration of the consonants but also in the degree of diphthongisation of the vowel.

The results described in the previous paragraph are compatible with a theory of the mental lexicon that postulates a direct connection between the context-specific meaning of a word token and the details of its phonetic realization. Such a theory has been computationally implemented as the Discriminative Lexicon Model, henceforth **DLM** (Baayen et al., 2019; Chuang and Baayen, 2021; Heitmeier et al., 2023c). In this model, lexis and morphology are acquired through a process of error-driven discriminative learning that allows for fine-grained alignments between low-level properties of form and low-level properties of meaning, both operationalized as high-dimensional numeric vectors. The model captures relationships between these vectors in two networks: a comprehension network that maps word form onto word meaning, and a production network that maps word meaning onto word form. The model does not store whole word representations of either kind; rather, the model’s memory consists of a set of networks in which the connection weights are continuously recalibrated with each learning event. In the corresponding theory of the mental lexicon, word forms and meanings do not have representations in memory. The forms are ephemeral auditory or visual experiences, which dynamically generate corresponding, equally ephemeral, meaning representations. Conversely, a meaning conceptualized by a speaker at a given point in time is dynamically transformed into ephemeral representations driving articulation. In other words, the DLM posits a lexicon in which lexical items are neither static nor discrete. Rather, the lexicon is taken to consist of a series of dynamic, modality specific neural networks (Baayen et al., 2019) which are constantly fine-tuned, by adjusting connection weights, in order to optimize word comprehension and production.

At this point, the question arises of how to understand the linguistic term ‘word’. In this study, we define **word token** as a pairing of a specific form (which can be represented mathematically as a high-dimensional numeric vector) with a specific meaning (which can also be represented mathematically as a high-dimensional numeric vector). We define **word type** (or simply, **word**) as a set of word tokens that have the property of having both similar forms and similar meanings. For instance, the set of phonetic realizations of 酒 *jiǔ* and their corresponding context-specific meanings (‘wine, liquor, spirits, alcoholic beverage’, <https://www.pleco.com/>, s.v.) jointly constitute the tokens of the word type 酒.² This working definition of ‘word’ seeks to do justice to the fact that no two tokens of the same word, as produced by humans, are ever completely identical in form. It also seeks to do justice to the insights from distributional semantics that what a word means varies with its context (Elman, 2009; Firth, 1968; Harris, 1954; Landauer and Dumais, 1997). Thus, in the framework of the DLM, words are sets of input-output pairs on which the production and comprehension networks are trained, that leave ‘traces’ in the connection weights of these networks, but that are themselves not stored as independent entities.

If the lexicon consists of networks of connection weights, then it is possible that properties of these networks can account for the variation in fine phonetic detail described above. This hypothesis is supported by a series of studies that have used the DLM to model such variation. For example, Gahl and Baayen (2024) used the DLM to model the different spoken word durations of English homophones, which turn out to be partly determined by how well the spoken form of a homophone token can be predicted from the meaning of that token. Using the same framework, Saito et al. (2023) showed that the strength of connection between features of form and meaning is predictive of the extent to which the tongue is lowered during the articulation of the German [a] vowel, as registered using electromagnetic articulography. For Mandarin, Chuang et al. (2023) were able to predict the duration of words in spontaneous speech in a similar way; as found for other languages, a Mandarin word’s duration tends to be longer if its form and meaning are well aligned, although this pattern appears to be restricted to shorter, semantically transparent words.³

Both the theoretical assumptions of the DLM and the empirical studies of English homophoneous words and affixes by Gahl and Baayen (2024) and Plag et al. (2017) respectively, suggest the possibility that Mandarin homophones also differ systematically in phonetic detail, i.e., that their segments and/or tones are realized slightly differently according to the intended meaning. This could apply to homographic pairs such as 大家 *dàjiā* ‘everyone’ and ‘art master’, as well as to non-homographic pairs such as 樹木 *shùmù* ‘tree’ and 數目 *shùmù* ‘number’. In other words, it is possible that the realizations of canonical tones are determined not only by the involuntary and voluntary constraints previously described, but also by the context-specific meanings of the word tokens on which they are realized. If this is correct, then conversely it is not only the four canonical pitch contours that help to distinguish between alternative meanings, but also the finer details of their phonetic realization. This brings us to the central hypothesis of our study:

²Note that, in this case, the Chinese orthographic character can be used to represent the word, since 酒 is not homographic: all its meanings cluster around the concept of ‘alcoholic beverage’.

³In this context, ‘shorter’ words are those whose canonical spoken forms have fewer segments.

The unique pitch contour of each spoken Mandarin word token is determined in part by the specific meaning of that token. A speaker can learn to produce meaning-specific pitch contours, and these meaning-specific contours are functional for the listener.

The following four predictions about Mandarin words follow from this hypothesis:

1. Variation in the tonal realization of a word cannot be reduced to the segment-related constraints on articulation previously described in the literature.
2. Information about a word’s meaning in context will improve prediction of its tonal realization.
3. Given a pitch contour, the meaning of its carrier token can be predicted above chance level, assuming the listener has previous experience of that word type.
4. Assuming they have previous experience of a given word type, a speaker can produce an appropriate pitch contour for a meaning they want to convey with that word.

In this paper we explore these predictions for disyllabic words with the canonical tone specification of a rising tone (T2) followed by a falling tone (T4), henceforth **RF**. The disyllabic word is a natural choice for our study since Mandarin vocabulary is composed mostly of disyllabic words (Huang et al., 2010; Wu et al., 2023). We decided to focus on the RF pattern, because it is the heterogeneous tonal combination with the highest number of word types and tokens in the speech corpus we used, namely the Taiwan Mandarin spontaneous speech corpus (Fon, 2004). We wanted to investigate a heterogeneous tonal combination rather than a homogeneous one to ensure that the results obtained are not specific to a given tone.

The remainder of the paper proceeds as follows. Section 2 addresses the first two predictions listed above. It describes how we used generalized additive modeling to analyze the pitch contours of RF words extracted from the above-mentioned corpus of spontaneous speech. We discuss our modeling strategy, and present the results of an analysis based on word type and one enhanced with word sense. These results are then triangulated with a separate analysis using a Random Forest algorithm. Section 3 addresses the third and fourth predictions. It describes how we used computational modeling with the DLM to demonstrate that meaning-specific pitch contours have the potential to facilitate comprehension and to be produced in response to intended meaning. Section 4 is a discussion of the implications of our results.

2 Establishing word and meaning-specific pitch contours

This section describes how we addressed the first two predictions outlined in Section 1. We first modeled the pitch contours of spoken tokens of Mandarin disyllabic words with the RF tonal pattern, using generalized additive modeling. To explore Prediction 1, we evaluated the effectiveness of word type as a predictor of tonal realization, compared with the segment-related articulatory constraints previously described in the literature. To explore Prediction 2, we evaluated whether adding information about a word’s meaning in context would improve prediction of tonal realization, compared with word type alone. Finally, we triangulated the results of the generalized additive models with a random forest analysis. Sections 2.1 to 2.4 describe aspects of the methodology, Sections 2.5 to 2.7 report the results of the generalized additive models, Section 2.8 presents the random forest analysis, and Section 2.9 summarizes the Section 2 overall.

2.1 Generalized additive modeling

Classical analyses of pitch typically take measurements at various contour landmarks, such as maximum and minimum F0 values. However, since pitch actually varies continuously with time, such analyses miss much of the detail of the F0 contour. To better capture the complete shapes of tonal variations, we modeled the pitch contours of the tokens in our data using generalized additive models, henceforth **GAMs** (Wood, 2017). GAMs relax the regression assumption that the relation between a predictor and response should be linear; instead, the model incorporates individual, potentially nonlinear relationships between each predictor variable and the response variable. For main effects, this relationship is estimated using functions known as smoothing splines (henceforth **smooths**), which can fit either a line or a (possibly wiggly) curve to the data, as required. Nonlinear interactions can be included using functions called **tensor product smooths**, which fit a wiggly (hyper)surface for the joint effect of two or more predictors. In addition, it is possible to include nonlinear random effects, for instance by using functions called **factor smooths** (Baayen et al., 2022), which fit a wiggly curve for each level of the random factor, e.g. for each individual speaker in the case of a by-speaker factor smooth. Because GAMs model complex non-linear relationships, they make it possible to model F0 as a nonlinear function of time across an utterance, while also including other predictors known to affect pitch, such as speech rate and speaker gender. They can thus capture fine-grained modulations of

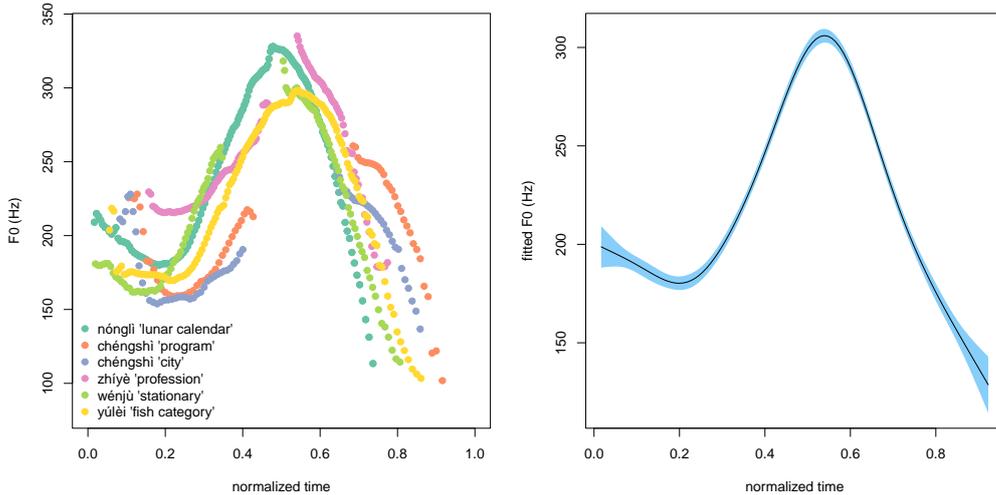


Figure 1: Toy dataset. The left-hand panel shows the F0 contours of single tokens of six Taiwan Mandarin words with the RF tonal pattern, produced in isolation by the same speaker. The right-hand panel shows the RF contour predicted by a simple GAM, using a thin plate regression spline smooth for normalized time as predictor.

pitch as time unfolds. For previous applications of GAMs to the analysis of pitch, see Chuang et al. (2021); Kösling et al. (2013); Sosluthy (2021); Sun and Shih (2021); Wieling et al. (2016).⁴

To illustrate our GAM-based modeling strategy, we created a toy dataset consisting of six disyllabic Mandarin word types with the RF tonal pattern. Their audio files were downloaded from *Meng Dian*, a publicly available Taiwanese online dictionary⁵. In the audio files, a single speaker pronounces each word twice, so that we had a total of 12 tokens all from the same speaker. We measured the F0 of each token every 15 ms and then transformed these time points onto a normalized time scale of 0-1. The F0 values of one of the two renditions of each of the six words are plotted on this normalized time scale in the left-hand panel of Figure 1. It can be seen that the RF tonal sequence in Mandarin is realized with a small initial fall, followed by a rise, and finally a much larger fall. The dipping realization of T2 is consistent with previous findings that the rising portion of Mandarin T2 is usually preceded by a slight fall (Ho, 1976; Moore and Jongman, 1997; Shen and Lin, 1991; Shih, 1988; Tseng, 1981; Xu, 1997).⁶

Using the **mgcv** package (Wood, 2017) for R (R Core Team, 2022), we fitted a GAM to the toy dataset, with F0 as the dependent variable and normalized time as the only predictor. Including time as a predictor allows us to model the entire pitch contour of a tonal pattern by predicting F0 at regular intervals across a token of that pattern. The pitch contour predicted by the GAM, shown in the right-hand panel of Figure 1, captures the general trend with some precision, mirroring the raw data on the left. This graph is the model’s best estimate of the average population contour for words with the RF tonal pattern, given the 12 tokens in our toy dataset. However, the empirical contours show considerable variation around this average, even for a single speaker producing citation forms in isolation. This variability in realization is the focus of our study.

To investigate whether individual tonal realizations might be related to the meaning of the carrier word token, we enriched the model with a by-word factor smooth, which is effectively a nonlinear, time-dependent, random effect for word type. In the present example, provided purely for illustration of the method, the by-word smooths are based on just two tokens of each word.⁷ This mixed model predicts, for each word type, a word-specific adjustment contour which has to be added to the population pitch contour to obtain the predicted pitch for a given word type. The by-word adjustment contours estimated by the GAM are visualized in the left-hand panel of Figure 2. The dotted line at $y = 0$ is a reference line: an adjustment curve for a given word that followed this line would indicate that no adjustment is needed and that this word’s pitch is identical to the population contour. Deviations above this reference line indicate an upward F0 adjustment, and deviations below it indicate a downward adjustment. The word 職業 *zhíyè* ‘profession’,

⁴Examples of the application of GAMs to other phonetic data such as tongue movement trajectories obtained from electromagnetic articulography can be found in Saito et al. (2023); Tomaschek et al. (2019b); Wieling et al. (2016). For GAMs used for modeling chronometric data, see, e.g., Baayen et al. (2022, 2017); Heitmeier et al. (2023c).

⁵Available at <https://racklin.github.io/moedict-desktop/addon.html>.

⁶The initial fall could also reflect dialectal variation specific to Taiwan Mandarin, as in this language, T2 is predominantly realized with a concave contour. This concave contour may have become a standardized realization that no longer reflects articulatory constraints (see e.g., Fon and Hsu, 2007).

⁷For detailed discussion of the ways in which these smooths can be specified, and the accuracy of these smooths, see Baayen et al. (2022).

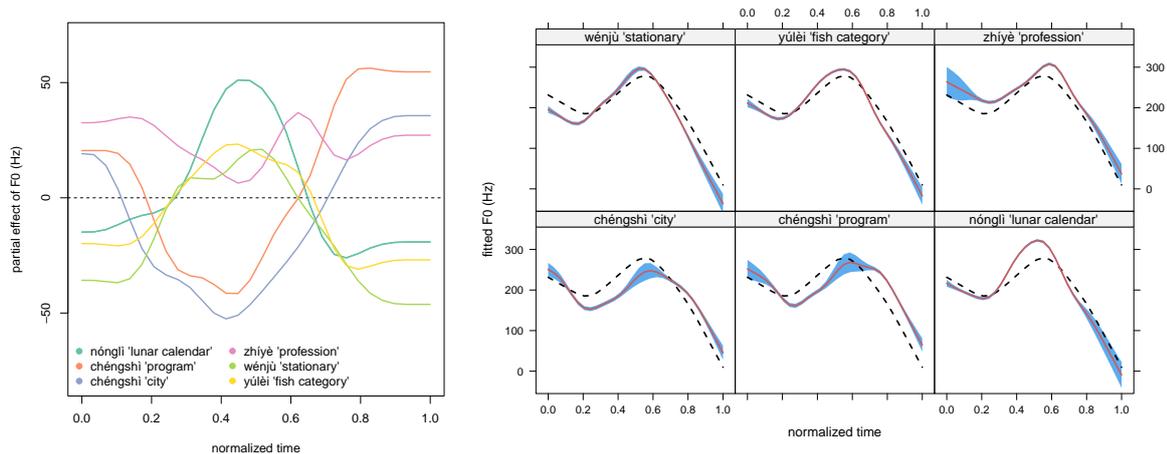


Figure 2: Toy dataset. The left-hand panel shows by-word adjustment contours from the toy model with only by-word factor smooth and normalized time as predictors. The right-hand panel plots the fitted contour for each word, with the predicted general contour (identical for all words) indicated by the dashed line.

for example, represented by a pink curve, requires an upward adjustment for the entire contour, although the amount of adjustment varies across time. When a given word’s adjustment contour is added to the general contour, we obtain its fitted contour, as shown in the right-hand panel of Figure 2. The dashed line in each graph represents the general contour which is, by definition, the same for all words. The red line (along with its confidence interval in blue) plots the fitted contour for the word in question. These fitted contours vary from word to word. For 職業 *zhíyè* ‘profession’, presented in the right-most upper panel, the entire fitted contour is above the general average contour, as expected. The homophone pair 程式 *chéngshì* ‘computer program’ and 城市 *chéngshì* ‘city’, shown in the left and middle lower panels, have similar but not identical fitted contours, as would be expected if word meaning, as well as word form, plays a role in determining tonal realization.

Although this is just a toy example, it does illustrate two important aspects of the more detailed analyses reported below. Firstly, we can decompose the observed pitch contour of any token into a general population contour plus various more specific contours, including a meaning-specific adjustment contour. Secondly, GAMs can identify such meaning-specific contours and, given an adequate sample size, they could inform us about whether including meaning-specific contours improves model fit.

In what follows, we turn to a much larger dataset of spontaneous spoken Taiwan Mandarin, and consider a much broader set of predictors that allow us to bring under statistical control a wide range of constraints known to co-determine the realization of pitch. If there is indeed a semantic component to the tonal realization of Mandarin disyllabic word tokens, then by-word factor smooths should be well-supported even when relevant control variables are taken into account.

2.2 Data

We used the Taiwan Mandarin spontaneous speech corpus (Fon, 2004), which consists of about 30 hours of recorded interviews with 55 native Taiwan Mandarin speakers aged between 20 and 60 years, 31 identified as female and 24 identified as male. This corpus contains 94,783 disyllabic word tokens, 11,482 of which have the RF tonal pattern⁸. These 11,482 tokens represent 707 orthographic word types. More than two-thirds of the tokens belong to one of five types: 然後 *ránhòu* ‘and then’, 時候 *shíhòu* ‘during which’, 不會 *búhuì* ‘do not’, 還是 *háishì* ‘still, or’ and 一樣 *yíyàng* ‘likewise’. In order to avoid model predictions being heavily biased towards these five high-frequency words, we randomly sampled 300 tokens for each of these five types for inclusion in our dataset.⁹ We furthermore excluded words with fewer than 20 occurrences from the dataset, in order to avoid overfitting to low-frequency words. As a consequence, our initial dataset comprised 9,496 tokens across 53 word types.

Subsequently, we extracted the sound files of these tokens and measured their F0 values using Praat (Boersma and Weenink, 2019). F0 values were taken every 15 ms, and therefore tokens of longer duration had more measurement datapoints than shorter tokens. In order to make sure every token had sufficient datapoints for model fitting, we excluded extremely short tokens with fewer than six datapoints, which

⁸These tokens also include tone sandhi words containing 一 *yī* and 不 *bù*, which have T2 realizations when followed by T4 syllables.

⁹The figure of 300 was an arbitrary choice at the upper end of the frequency range for the other words in the dataset.

constituted about 5% of the data. Next, we removed tokens where an F0 extraction error was likely. F0 extraction errors usually result from pitch halving or doubling, and lead to abrupt big changes in the recorded F0 values. We therefore first obtained, for each token, all the F0 differences between consecutive measurements, and then calculated the standard deviation of these difference values. The standard deviation is large when F0 measurements are discontinuous and fluctuate abruptly. Tokens with standard deviations greater than the 9th decile of the distribution were considered to be outliers, hence likely to involve extraction errors, and were excluded from further analyses. Finally, two words were excluded because their tokens were all contributed by only one speaker.

The final dataset for the first analysis, reported below in Section 2.6, consisted of a total of 3,778 tokens representing 51 word types. Since these types do not include any heterographic homophones, there is a one-to-one correspondence between the orthographic labels of the tokens in our data and their canonical spoken forms. We therefore assume that tokens with the same label bear some similarity to one another in both form and meaning, i.e. belong to the same word type.¹⁰

2.3 Predictors

The core predictors in our GAM models are described in Section 2.3.1 below. As far as possible, we also included as controls all the variables that have previously been shown to influence tonal realization, as outlined in Section 1 above. These control predictors can be grouped into three major categories: speaker-related, context-related, and segment-related. They are described in Sections 2.3.2 to 2.3.4 respectively.

2.3.1 Core predictors

Word type: We coded each word token in our dataset for its word type (**word**), using the orthographic representation of the token in the corpus as the identifier of its word type.¹¹

Sense: Unlike heterographic homophony, homographic homophony and polysemy are common in Mandarin disyllabic words. In lexicography, such diversity of meaning is usually addressed by attempting to enumerate the various possible senses of a given orthographic form. Similarly, in computational semantics, systems have been devised for disambiguating word senses from amongst a finite set of possibilities. The validity of this approach has been questioned, e.g. by Kilgarriff (2007, p. 29), who pointed out that there are ‘no decisive ways of identifying where one sense of a word ends and the next begins’; polysemy is actually much more subtle and nuanced than a set of discrete possibilities would suggest. Nevertheless, sense annotations do capture, however crudely, some aspects of the variability in words’ meanings. Furthermore, within the context of modeling pitch contours with GAMs, discrete senses are convenient because we can straightforwardly estimate specific pitch contours for each sense. We therefore coded every word token in the dataset for **sense**, using a word sense disambiguation system (Hsieh and Tseng, 2020) based on the Chinese Wordnet (Huang et al., 2010). The possible values of this variable correspond to the senses identified by the disambiguation system. More than one sense was identified for 35 of the 51 words in the dataset, with a total of 130 senses overall. All except two of the words had between one and five senses; of the two outliers, one had 6 senses and the other had 9 senses. Note that, because the sense labels in our data are nested under the orthographic form, and there are no synonyms, **sense** includes all the information in **word**, plus additional information about the meaning of any given token.

Normalized time: The points in time at which F0 measurements were taken (at 15 ms intervals) were, for each token, transformed into a normalized time scale of 0-1 to produce the variable **time**.

2.3.2 Speaker-related controls

Gender: Speakers identified as female usually have a higher pitch register and wider pitch range than speakers identified as male. Furthermore, with respect to tonal realizations in Taiwan Mandarin, a number of studies have documented detailed gender-dependent differences in various sociolinguistic domains (Fu, 1999; Huang, 2008; Wu, 2009, 2003). We therefore included **gender** as a simple control variable to account for intrinsic pitch height and range differences between speakers of different genders, as labeled in the corpus, and also allowed **gender** to interact with **time**, to accommodate possible gender-specific modulations of the pitch contour.

¹⁰Heterographic homophones are relatively rare amongst Mandarin disyllabic words. By definition, the spoken forms of disyllabic words have more possible combinations of segments and canonical tones than monosyllabic words do; this means there is less need to reuse the same canonical spoken forms for different meanings.

¹¹Note that, although word frequency is an important predictor for response variables such as spoken word duration, it is not a factor that has been widely reported to co-determine the shape of pitch contours in Mandarin (but see Bi et al., 2015). We verified for our data that when word frequency is included in a model that also has access to word type, frequency is not significant, whereas word type is well-supported. In this study, frequency of use is therefore not discussed any further.

Speaker identity: Speaker identity was included to account for any idiosyncratic tonal realizations specific to individual speakers. We included `speaker` not only as a main effect, but also in interaction with `time`, using by-speaker factor smooths.

2.3.3 Context-related controls

Speech rate: The shape of tonal contours is modulated by how fast speakers talk. Generally, the faster the speech rate, the more likely it is that tonal contours will deviate from their underlying shapes (Cheng and Xu, 2015; Tang and Li, 2020). We therefore calculated `speech_rate` for each token, defined as syllables per second over the time period from four words before to four words after the target token, inclusive.

Adjacent tones: When a tone is expected to start at a different pitch from where the previous one ends, e.g., a falling tone followed by a high level tone, the degree of co-articulation, and hence deviation from the canonical tonal shapes, will be greater than when two tones are contiguous, e.g., a high level tone followed by a falling tone (Shih, 1988; Xu, 1994). In addition, although the details differ across studies, tonal co-articulation is usually found to be bi-directional, i.e., both anticipatory and preservatory (Huang and Chiu, 2023; Shen, 1990b; Xu, 1997). For our analyses, we therefore coded the tonal category of each token’s preceding and following syllables in the corpus. When a target token occurred utterance-initially or utterance-finally, the preceding or following tonal category was coded as ‘null’. This gave us six possible tonal categories for both the preceding syllable and the following syllable: four lexical tones, one neutral tone, and ‘null’. We therefore created the factor `adjacent_tone` with 36 levels to account for each possible combination of properties of the preceding and following syllables.

Utterance position: The realization of tone in an utterance is also influenced by sentence intonation (Ho, 1976; Shen, 1989, 1990a; Tseng, 1981). For example, statement intonation is often characterized by a downward trend, resulting in pitch declination (Shih, 1997). Question intonation, on the other hand, can potentially lead to a final rise, although this largely depends on the syntactic structure and/or emotive force of the question concerned (Chuang et al., 2007; Lee, 2005). For the current study, we simply calculated the normalized position of each token in the relevant utterance. We defined an utterance as a sequence of words preceded and followed by a perceivable pause (regardless of duration), as indicated by the labels provided in the corpus. The variable `utterance_position` is the position at which a given token occurs in an utterance divided by the total number of words in that utterance. This predictor is therefore bounded between 0 and 1. For utterances with only one word, the utterance position was set to 1.

Bigram probability: Bigram probability is a measure of a word’s contextual predictability based on its relative frequency of co-occurrence with the other words in its context; the higher the bigram probability, the more predictable a target word is in the given context. It has been found that a word’s phonetic realizations are intimately related to its contextual predictability. In general, higher predictability is associated with shorter word duration and a greater degree of spectral reduction (Bell et al., 2003; Gahl et al., 2012). Specifically for tonal realizations in Mandarin, there is some evidence that these too are sensitive to contextual predictability; when a word is more contextually predictable or represents given information, its F0 excursion and range are found to be diminished (Hsieh, 2013; Ouyang and Kaiser, 2015). In the present study, following Gahl et al. (2012), we calculated the bigram probabilities of target tokens in two ways: `bigram_previous`, based on the preceding word, and `bigram_following`, based on the following word. These two variables are defined, respectively, as follows:

$$P(w_n|w_{n-1}) = \text{Freq}(w_{n-1}, w_n) / \text{Freq}(w_{n-1}),$$

$$P(w_n|w_{n+1}) = \text{Freq}(w_n, w_{n+1}) / \text{Freq}(w_{n+1}),$$

where `Freq` denotes word frequency in the corpus of Taiwan Mandarin (Fon, 2004).

2.3.4 Segment-related controls

Vowel height: It has long been recognized that different vowels have different intrinsic pitch, a finding established for a great number of different languages, including Mandarin (Ho, 1976; Ladd and Silverman, 1984; Shi and Zhang, 1987; Whalen and Levitt, 1995). Specifically, high vowels tend to have higher F0 values than low vowels. For our disyllabic words, we coded the vowel heights of the vowels of the first and second syllables as two separate predictors, `vowel1` and `vowel2` respectively. For monophthongs such as /i/ and /a/, we distinguished between three vowel heights: ‘high’, ‘mid’, and ‘low’. For diphthongs such as /ai/ and /ei/, which are characterized by within-vowel changes in height, we added two additional levels: ‘low-high’ and ‘mid-high’. This means that there are theoretically 25 possible combinations of `vowel1` and `vowel2`. Our dataset included 20 of these possible combinations.

Onset: The effect of onset consonant on F0 has been studied in considerable detail in Mandarin. Ho (1976), for example, found that voiced onsets lead to lower F0 in the following vowel, as compared to voiceless onsets. Aspiration, on the other hand, often causes pitch perturbation, resulting in higher F0, although the magnitude of this effect appears to be tone dependent (Xu and Xu, 2003). Following Howie (1974), we distinguished onset types according to manner of articulation, voicing, and aspiration. For each of the two syllables in our target words, we distinguished between ‘aspirated-affricate’, ‘aspirated-stop’, ‘unaspirated-affricate’, ‘unaspirated-stop’, ‘voiceless-fricative’, and ‘voiced’. Syllables that do not have onsets and start with vowels or glides instead were coded as ‘null’. Our dataset contained 30 different combinations of the onset type of the first syllable (`onset1`) and that of the second syllable (`onset2`).

Rhyme structure: Although effects appear to be unstable and are not always reliably observed, some studies have reported variation in F0 for different Mandarin syllable types (Fon and Hsu, 2007; Howie, 1974; Xu, 1998). In our models, we therefore included a control variable for syllable structure. Given the strict phonotactic constraints governing the syllables of Mandarin, a syllable can maximally be composed of an onset consonant, a prenuclear glide, a nucleus vowel, and finally a coda consonant (Duanmu, 2007). In some theoretical descriptions, a coda consonant must be a nasal; in other descriptions, it can be either a nasal or a postnuclear glide as in the case of /ai/, for example. In this study, we coded the latter cases as diphthongs in the vowel height predictors, `vowel1` and `vowel2`, and only coded for a coda consonant when the syllable included a final nasal. For each of the two syllables in our target words, we therefore coded the structure of the rhyme as ‘V’, ‘GV’, ‘VN’, or ‘GVN’. This coding specifies, for a given syllable, whether there is a prenuclear glide, as well as whether there is a final nasal. Applied to the two syllables of our target words separately (`syllable1` and `syllable2`), we obtained 14 attested combinations of rhyme structures.

2.4 Modeling Strategy

Because pitch typically changes continuously and gradually across the time course of an utterance, it is inevitable that our response variable of F0 measurements is characterized by significant autocorrelation. That is, the F0 at time t is correlated with and can thus to some extent be predicted from the F0 at $t - 1$. Autocorrelation is particularly problematic for regression modeling because the residuals of an autocorrelated response variable are also unavoidably autocorrelated. This means that a central assumption of regression modeling is violated, namely that the residuals should be independent of one another.

In GAMs, the issue of autocorrelation can be addressed by incorporating a first order autoregression model for the errors, denoted as **AR(1)**. An AR(1) model is a linear model that predicts a given value of a time series from the immediately prior value; including an AR(1) process in a GAM enables the model to accommodate structure in the residuals by positing a linear relationship between a given residual and its preceding residual. This can be a highly effective way of dealing with autocorrelation. However, for the current dataset, including an AR(1) process led to two new problems (see Appendix A for further details). Firstly, the AR(1) process in a GAM assumes that the degree of autocorrelation is invariant; however, the degree of autocorrelation actually varies considerably across the tokens in our dataset and the AR(1) process therefore over-corrects for some tokens but under-corrects for others. Secondly, the AR(1) process has an undesirable effect on the overall distribution of the residuals in our models, increasing the extent to which they deviate from normality; this happens because the data contains a large number of discontinuous contours arising from voicelessness or creaky voice (see Figures A3 and A4). For these two reasons, we decided not to incorporate an AR(1) process in our modeling.

The decision not to include an AR(1) process means that we have fewer independent datapoints than a GAM assumes. Our models are therefore anti-conservative and calculate confidence intervals that are too narrow. For this reason, in the analyses that follow, we do not report p-values for the partial effects in our GAMs. Instead, we adopted a two-fold modeling strategy. Firstly, in order to evaluate whether a given predictor is relevant for understanding F0 contours, we made use of Akaike’s Information Criterion (AIC), assessing the extent to which AIC decreases, i.e. model fit improves, when a predictor is added to a baseline model. Secondly, we investigated the adequacy of our predictors by means of cross validation. That is, we held out a small portion of our dataset as testing data, fitted our models to the remaining data, and then assessed model accuracy on the testing data. In this way, we can assess the precision with which our models predict novel, previously unseen, data, and establish whether inclusion of a predictor improves prediction accuracy.

We used this modeling strategy to explore the first two predictions presented in Section 1, repeated here for convenience:

1. Variation in the tonal realization of a word cannot be reduced to the segment-related constraints on articulation previously described in the literature.
2. Information about a word’s meaning in context will improve prediction of its tonal realization.

Prediction 1 is based on the hypothesis that the unique pitch contour of each spoken Mandarin word token is determined in part by the meaning of that token. If this is correct, then variations in tonal realization arise not only from the mechanical constraints on articulation previously described in the literature, but also as a result of form-meaning connections in the lexicon. We therefore investigated whether using `word` as a predictor would lead to a more precise model of tonal realization, as compared to a model using the set of segment-related predictors (e.g., vowel height) that have previously been identified as relevant. Since `word` incorporates not only the segmental form of a word type but also the associated semantics, our expectation was that it would be a superior predictor compared to all the previously identified segment-related variables considered jointly, when we controlled for contextual effects such as speech rate.

As a first step, we fitted a baseline GAM to log-transformed F0, including all the speaker-related and context-related control variables described in Sections 2.3.2 and 2.3.3, in interaction with `time`. We then fitted six further models, each with one segment-related control variable added to the baseline, allowing us to investigate the articulatory effects described in the literature. To compare these effects jointly against the effect of `word`, we fitted two additional models, one with all six segment-related control variables included in addition to the baseline, and the other with only `word` as an additional predictor. Both `word` and the segment-related controls (`vowel1`, `vowel2`, `onset1`, `onset2`, `syllable1`, `syllable2`) were modeled as factor smooths in interaction with `time`.

Prediction 2 is based on the hypothesis that variation in tonal realization serves to make word tokens with different meanings more discriminable from one another. To address this prediction, we compared the model with `word` as the sole predictor added to the baseline against a model in which `sense` was the sole predictor added to the baseline. Our expectation was that if semantics is indeed at issue, replacing `word` by `sense` should improve model fit even further. In the dataset for the word-type analysis, the frequency distribution of `sense` is skewed towards the right, with about half of the senses having no more than 13 tokens. To make sure that all senses included in the models had sufficient tokens for statistical evaluation, we therefore used only a subset of the data for the models used to investigate the effect of `sense`. Since the median number of tokens per sense in the dataset was 13.5, we excluded senses with fewer than 14 tokens. This left us with a dataset of 3,458 tokens representing 65 senses across a smaller set of 48 word types. We used this smaller dataset for models evaluating `sense` as a predictor.¹² The statistical analysis proceeded in the same way as described above for word type, except that we added the additional model with `sense` as the only predictor in addition to the baseline.

We fitted our models using the `mgcv` package (Wood, 2017) in R (R Core Team, 2022). The distributions of the residuals of Gaussian models fitted to the pitch contours were t-distributed rather than normally distributed, and we therefore set the `family` argument to `scat` (scaled-t, Wood et al., 2016).

2.5 The baseline GAM

We fitted the baseline GAM to log-transformed F0, including all speaker-related and context-related predictors, using the following model specification:

```
pitch ~ gender + s(time, by = gender) +
  s(time, speaker, bs = 'fs', m = 1) +
  s(speech_rate) + ti(time, speech_rate) +
  s(utterance_position) + ti(time, utterance_position) +
  s(bigram_previous) + ti(time, bigram_previous) +
  s(bigram_following) + ti(time, bigram_following) +
  s(time, adjacent_tone, bs = 'fs', m = 1)
```

The first line of this model requests a main effect for gender, in order to account for male voices being lower on average than female voices. In addition, the model requests a separate smooth for each gender. The upper left-hand panel of Figure 3 plots the predicted contours for speakers identified as female (red) and those identified as male (blue). Similar to the pattern observed in read speech (cf. Figure 1), the realization of the RF tonal pattern in spontaneous speech is characterized by a shallow fall, followed by a long rise, and finally a much larger fall. Speakers identified as male show reduced pitch excursion compared to those identified as female, presumably due to male voices having a more compressed pitch range. The second line of the model requests by-speaker nonlinear random effects, using factor smooths. These factor smooths specify, for each speaker, the specific way in which that particular speaker modulates the general F0 contour associated with their gender.¹³

¹²For completeness, we note that 35 words in the smaller dataset have only one sense. Of these 35 words, six have only one sense listed in the Chinese Wordnet, and seven are not included in the vocabulary of the Chinese WordNet. For the rest, either the tagger only identified one sense, or only one sense of the word had more than 13 tokens in our dataset.

¹³The syntax `bs='fs'` on the second line of the baseline model specification has a similar effect to the `by` argument on the first line; both terms request a smooth for each level of a single factor variable (in this case, `speaker` and `gender` respectively).

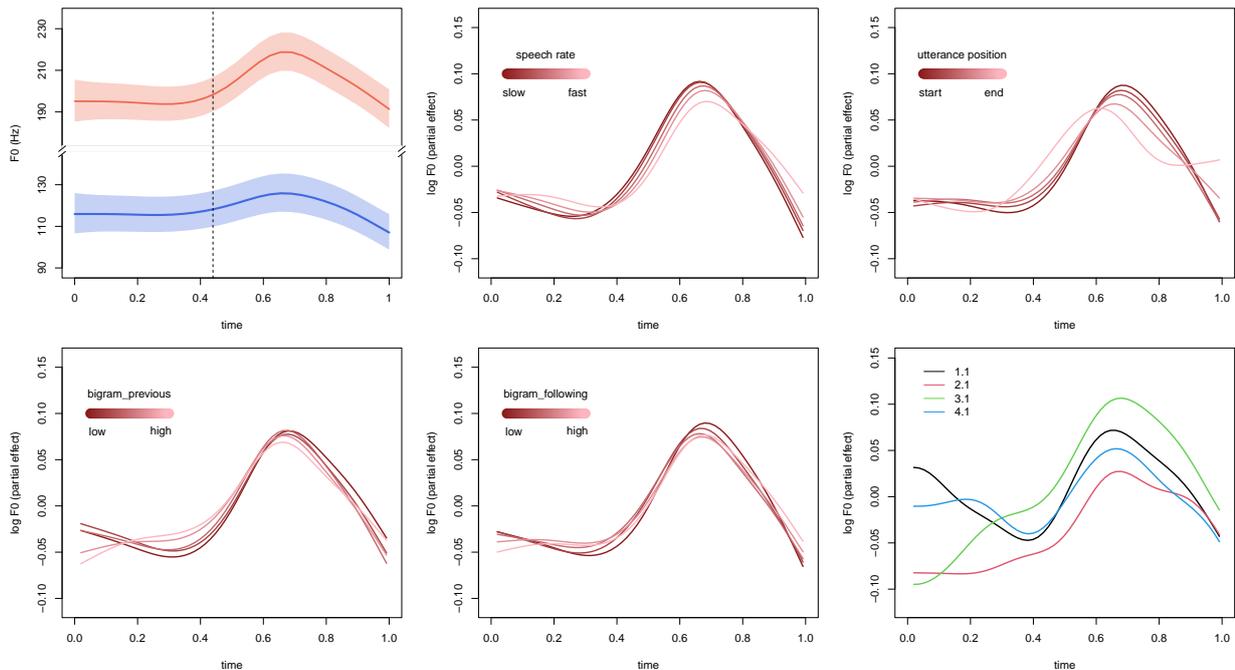


Figure 3: Baseline GAM. The upper left-hand panel shows the predicted base contours for speakers identified as female (red) and speakers identified as male (blue). The next four panels show, for female speakers, how the base contour is modulated by speech rate, utterance position, previous bigram probability, and following bigram probability, respectively. The final panel presents, again for female speakers, the effect of tonal coarticulation with the tone of the preceding word, given that the following word has a high-level tone.

The next four lines in the model formula deal with the four numerical context-related controls, namely `speech_rate`, `utterance_position`, `bigram_previous` and `bigram_following`. For each of these variables, the model requests a main effect smooth in combination with a tensor product smooth for the interaction of the given variable with `time` (using an interaction-specific tensor product smooth specified with `ti`). The upper mid panel of Figure 3 plots the modulating effect of speech rate on the base contour of speakers identified as female. Higher speech rates, represented by lighter shades of red, reduce the amplitude of the wave. The effect of position in the utterance is depicted in the upper right-hand panel of Figure 3. The tonal shape is clearly most different when the word occurs towards the end of an utterance, in which case we observe an earlier peak. This might be due to the fact that words in singleton utterances are treated as occurring at the end of an utterance. The left and mid panels in the lower row of Figure 3 present the effects of the bigram probabilities given the preceding and following word respectively. Similar to the effect of speech rate, higher bigram probabilities, represented by lighter shades of red, generally lead to a smaller F0 excursion.

The final line of the model specification requests factor smooths for `adjacent_tone`, requesting a separate smooth for each of its 36 levels. The effect of `adjacent_tone` is presented in the lower right-hand panel of Figure 3 for those tokens which have T1 as following tone. Unsurprisingly, the four predicted contours end similarly. However, the initial part of the contour diverges considerably, depending on the preceding tone. As expected, tonal context has a very large effect on the shape of the F0 contour.

In what follows, we take this model as our baseline model, with all contextual covariates controlled for, and compare the effects of predictors representing individual aspects of word form with the effect of word type as a whole.

2.6 Results and discussion: word type

2.6.1 Evaluation of predictors

The left-hand panel of Figure 4 presents the improvement in model fit, as compared to the baseline model, for the six models with a single additional segment-related control, the model with all segment-related

The `by` argument is generally preferred when the factor has few levels, and the levels are of interest, as in the case of `gender`. For computational efficiency, the `fs` argument is preferred when there are many levels and these levels are of less direct interest, as in the case of `speaker`. When specifying a `by`-smooth, a separate term requesting a main effect for the intercept needs to be specified. In contrast, a factor smooth incorporates adjustments to the intercept, thus effectively calibrating the individual smooths for their relative position with respect to the general intercept.

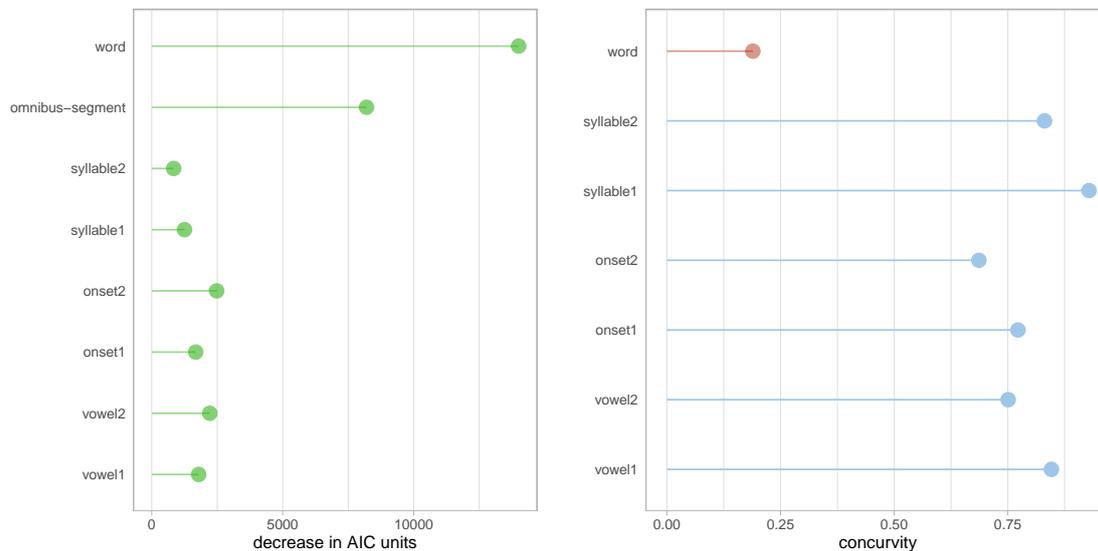


Figure 4: Word-type analysis. The left-hand panel shows model fit improvement, gauged by decrease in AIC units when a predictor (or set of predictors) is added to the baseline model. The right-hand panel shows the concurvity score of individual predictors in two models: the omnibus-segment model with all segment-related control variables (blue) and the word-type model with **word** as predictor (red).

controls (henceforth the **omnibus-segment** model), and the model with only **word** as additional predictor. Improvement in model fit is gauged by the magnitude of any decrease in AIC. As can be seen, all segment-related controls improve model fit, a result that dovetails well with previous studies.

The inclusion of all six segment-related controls jointly in the omnibus-segment model provides a substantially better fit than any single predictor. However, adding only **word** to the baseline model results in an even better fit, the difference in AIC being no less than 5804 AIC units. Although the segment-related controls address all the word properties previously found to influence tonal realization, the contribution of just **word** by itself to the model fit is much stronger. It is clear that the association between word type and tonal realization cannot be reduced to the segment-related constraints on articulation previously described in the literature. The actual pitch contour is richer than what can be predicted from all phonetic features that have been found to be relevant.

The omnibus-segment model has the problem that its key predictors are correlated with one another. In a GAM, the nonlinear equivalent of collinearity is **concurvity**. The concurvity score of a predictor is a number bounded between zero and one that measures the degree to which the effect of a given independent variable can be predicted by one or more of the other independent variables in the model. If a predictor's concurvity is low, this predictor has its own explanatory value. However, if the concurvity is high, the predictor is strongly confounded with other predictors.

The right-hand panel of Figure 4 presents the concurvity scores of the segment-related controls in the omnibus-segment model (marked in blue) and the concurvity score of the predictor **word** in the word-type model (red). For the omnibus-segment model, concurvity scores of all predictors are high, which is perhaps unsurprising given the restrictive phonotactic constraints governing Mandarin syllables (Duanmu, 2007). The high concurvity indicates that the effects of the segment-related controls are confounded with one another, rendering interpretation of the individual effects difficult if not impossible. On the other hand, in the word-type model, the concurvity of **word** is low, so that the interpretation of the effect of individual word types on the F0 contours is straightforward.

The predicted pitch contours of a sample of 15 word types are presented in Figure 5. To better visualize how the word-specific tonal modulations differ from one another, the partial effect predicted for each word has been added to the general contour for speakers identified as female (cf. Figure 3). In general, the fall-rise-fall pattern can be observed for all these words, but the details of tonal excursions differ significantly from word to word. For example, while the initial falling part is very prominent for words like (c) 決定 *juédìng* 'decision' and (f) 全部, *quánbù* 'all', it is rather muted for (m) 一半 *yībàn* 'half' and (d) 年紀 *niánjì* 'age'. In addition, in terms of the degree of undulation, some words have more reduced tonal range, such as (a) 不是 *búshì* 'not' as compared to (j) 文化 *wénhuà* 'culture', which has an extensive F0 excursion.

A closer inspection of Figure 5 reveals that, as expected, some of the word-specific contours appear to be consistent with the words' canonical segmental properties. For example, the initial fall appears to be more salient when the onset of the first syllable is a voiceless sibilant, e.g., (k) 習慣 *xíguàn* 'habit' or an

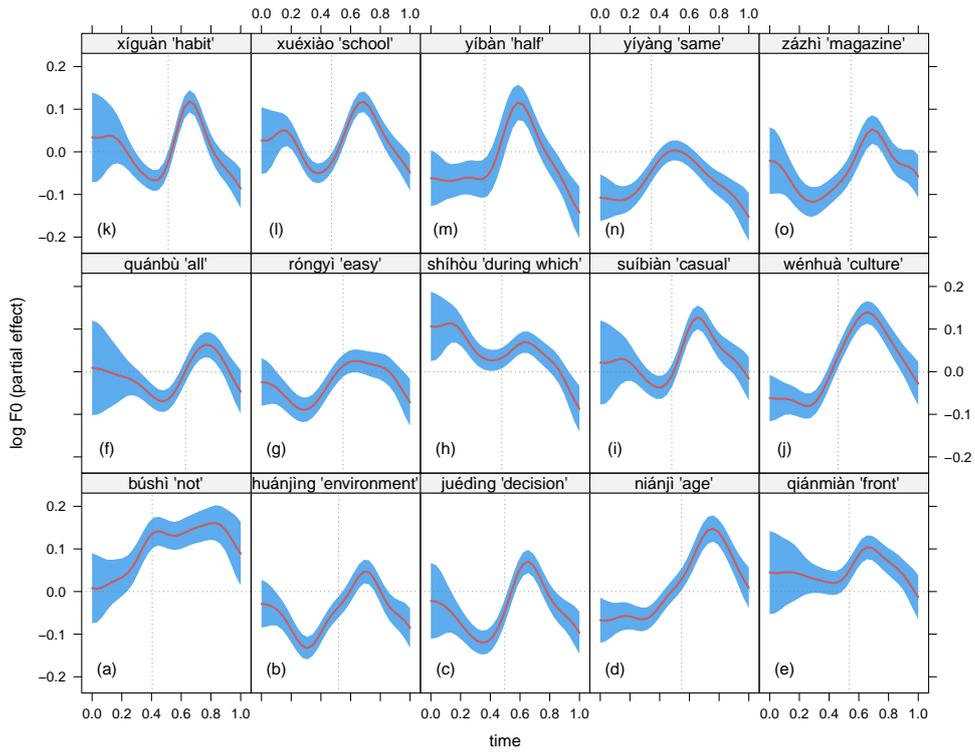


Figure 5: Word-type analysis. Examples of the pitch contours predicted by the general smooth for time for female speakers, combined with the partial effects of the factor smooth for *word*. These partial effects do not include the general intercept, nor the differences in pitch between female and male speakers. As they represent the pure effect of *word* on the pitch contour, irrespective of other predictors, the curves are centered around the y-axis (indicated by a horizontal dotted line). The vertical dotted lines in the panels indicate the average (word-specific) syllable boundary.

affricate, e.g., (e) 前面 *qiánmiàn* ‘front’, with the onsets / ϵ / and / $t\epsilon^h$ / respectively. However, it should be kept in mind that the model is imputing F0 values for these voiceless onsets, where no periodic wave form is actually produced. In Figure 5, the plotted 95% confidence intervals for the early timesteps in these words can be seen to straddle the horizontal axis, despite the fact that the confidence intervals are anti-conservative and therefore appear narrower than they should (cf. Section 2.4). Since the horizontal axis represents no effect, there is no good evidence for modulation of the general F0 contour early on in these words. Another segmental property that is to some extent visible in Figure 5 is the length of the second syllable, relative to the length of the first syllable. If the second syllable is relatively short as compared to the first syllable, e.g., (f) 全部 *quánbù* ‘all’ and (g) 容易 *róngyì* ‘easy’, the final fall tends to be attenuated, as expected given that relatively less time is available to physically implement a large fall in pitch. Nevertheless, the superior performance of **word** over the segment-related controls in our models indicates that such articulatory effects are not the only elements at play in determining tonal realization.

2.6.2 Cross validation

We have seen that word type is an excellent predictor of pitch realization, outperforming the segment-related predictors previously identified in the literature. However, due to the presence of strong temporal autocorrelations inherent in the pitch measurements, the confidence intervals in our models, shown for example in Figure 5, are anti-conservative and are therefore far from optimal for estimating the uncertainty of the predicted partial effects. Nevertheless, the models do make explicit predictions of F0 contours, and can do this even for held-out data, i.e. tokens that were withheld from the model during model fitting (training). If our hypothesis is correct, a GAM that has access to **word** should provide superior prediction accuracy on held-out data compared to a GAM that has access only to the segment-related controls. We therefore evaluated prediction accuracy under cross validation. We first held out 10% of the current data as test data, and used the remaining 90% as training data. We ensured that every word type was represented in both the training and test data, and that the number of tokens per type in the test data was proportional to that in the training data.

We fitted nine models to the training data. In addition to the eight models already introduced above, and assessed in Figure 4 (left-hand panel), we added one more model that was given data in which the values of **word** were randomly permuted. That is, tokens of a given word were now assigned different random word labels. In what follows, we refer to this model as the **random-word** model. If the effect of **word** is genuine, then random permutation of the word labels should significantly reduce prediction accuracy. To quantify model accuracy, we obtained the models’ predictions for the F0 contours of the held-out test data, and calculated the sum of squared errors (SSE) as measure of prediction accuracy.¹⁴ The SSE for the held-out data of a given model should be smaller than the SSE of the baseline model if the addition of one or more predictors indeed improves that model’s prediction accuracy. We implemented cross-validation using both the scaled t-distribution (**scat**) and simple Gaussian models. As models with **scat** are computationally expensive to estimate, for environmental reasons we used these models for one held-out dataset only and compared the results with those of less expensive Gaussian models for the same held-out data. Since prediction accuracy turned out to pattern similarly for **scat** and standard Gaussian models, and since fitting Gaussian models requires much less computational resource, we ran Gaussian models to cross-validate our results for 100 random held-out datasets.

Cross-validation results are presented in Figure 6. The left-hand panel presents the SSE difference between the baseline model and each of the nine models of interest, for one random held-out dataset. Larger positive values indicate that the model of interest offers more precise predictions than the baseline. The difference in SSE obtained with a standard Gaussian model is shown in blue, and the difference obtained with a **scat** model is shown in orange. All the individual segment-related controls increase prediction accuracy over the baseline, albeit to varying degrees. However, the omnibus-segment model and the word model produce significantly greater increases in prediction accuracy, with the latter reducing the SSE to a larger extent than the former, replicating the model fit results. Moreover, when word labels are randomized, model accuracy plummets: the SSE of the random-word model is greater than that of the baseline model.

Because the results of GAM models with and without **scat** are very similar for the first set of test data, we infer that the distribution of results across the 100 cross-validation runs using Gaussian models will also be informative. These results are reported in the right-hand panel of Figure 6. The overall pattern is very similar to the result of the initial single cross-validation run. The word model is still far more accurate than all the other models, and significantly outperforms the omnibus-segment model. In other words, the greater reduction in AIC that we observed for the word model in Section 2.6.1 is not just a side-effect of an evaluation procedure that happens to favor that model over the omnibus-segment model. The word model truly provides enhanced predictions.

¹⁴The sum of squared errors (SSE) is the sum of the squared difference between the observed and predicted values. A smaller SSE indicates more precise model predictions.

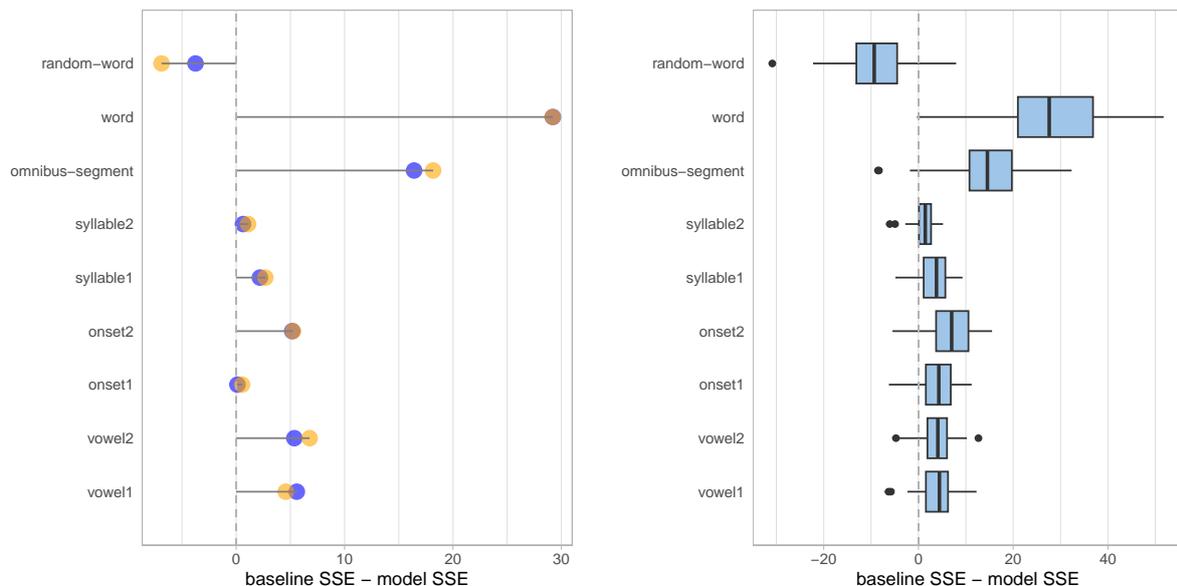


Figure 6: Model accuracy under cross validation for the word analysis. The left-hand panel presents the reduction in the sum of squared errors (SSE) for models using one set of test data. Positive values indicate lower SSE as compared to the SSE of the baseline model. Orange dots represent scat models, while blue dots represent Gaussian models. The right-hand panel shows the results of 100 runs of cross validation, using Gaussian models. The boxplots represent the distributions of reduction in SSE.

The results presented so far provide strong evidence that word type is predictive of tonal realization, over and above the segment-related predictors established by previous studies. Our hypothesis, arising from the theoretical framework of the Discriminative Lexicon Model (DLM Baayen et al., 2019; Chuang and Baayen, 2021; Heitmeier et al., 2023c), is that the predictive power of word type arises not only from articulatory constraints, but also from a close association between word meaning and phonetic form, which enables the language learner or user to discriminate more efficiently between forms with different meanings. However, since word meaning varies with context, there is no one-to-one correspondence between word type and the meaning of a given token; **word** can only encompass a rather general approximation of what each token means. If word meaning is indeed predictive of tonal realization, then a model replacing **word** by **sense** should improve model fit even further.¹⁵

2.7 Results and discussion: sense

2.7.1 Evaluation of predictors

Figure 7 shows model fit improvement and concurrency for models based on the smaller dataset that included at least 14 tokens of each word sense. Even with this smaller dataset, the overall pattern of results is very similar to that of the word type analysis. Critically, however, **sense** appears to be a somewhat better predictor than **word**. Not only does it account for more variance in F0, resulting in better model fit, but its effect is also less confounded with other covariates in the model, as indicated by a smaller concurrency score.

Figure 8 presents the predicted tonal contours for different senses of three words: 不要 *búyào* (left), 實在 *shízài* (upper right), and 能夠 *nénggòu* (lower right). The word 不要 *búyào* is a polysemous negation marker in Mandarin. The four senses that are found in our dataset are ‘prohibition’, ‘dissuasion’, ‘unnecessity’, and ‘to wish something to not happen’ (s1 to s4, respectively). It can be seen that the different senses have clearly different tonal realizations. The panels on the right-hand side of Figure 8 present the predicted contours for the other two words, each of which has two senses in our data. For 實在 *shízài*, tonal realizations vary greatly between the two senses (‘truly’ and ‘indeed’), whereas the realizations of the two senses of 能夠 *nénggòu* (‘being capable’ and ‘enabling’) are more alike, and differ mainly with respect to the amplitude of the pitch inflection.

¹⁵Note, however, that as discussed in Section 2.3, the senses that constitute the possible values of our **sense** variable discretize a much more subtle and interesting palette of shades of meanings.

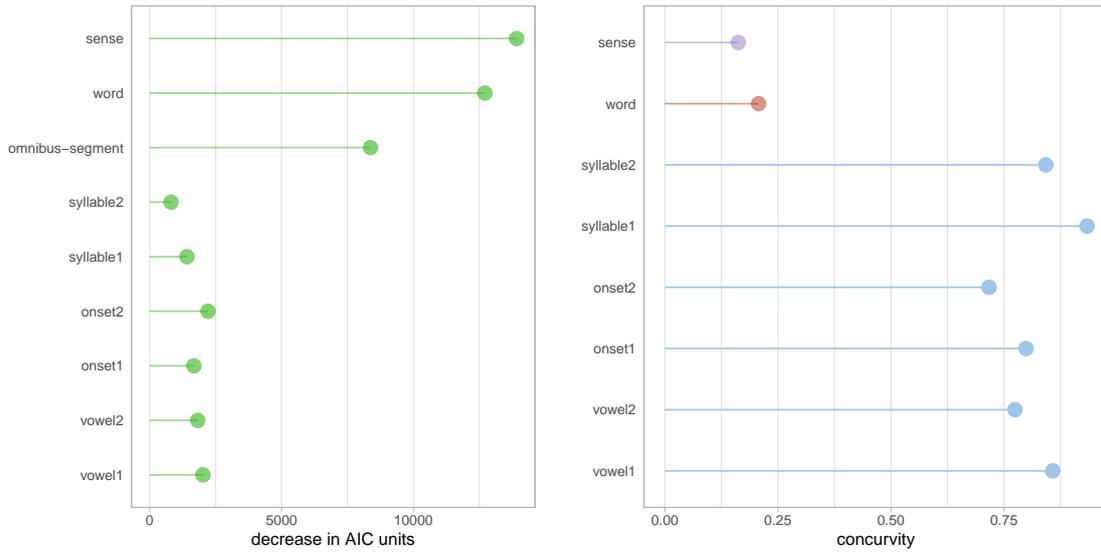


Figure 7: Sense analysis. The left-hand panel shows model fit improvement, gauged by decrease in AIC units when a predictor (or set of predictors) is added to the baseline model. The right-hand panel shows the concurvity score of individual predictors in three models: the omnibus-segment model with all segment-related control variables (blue), the word-type model with *word* as predictor (red), and the sense model with *sense* as predictor (purple).

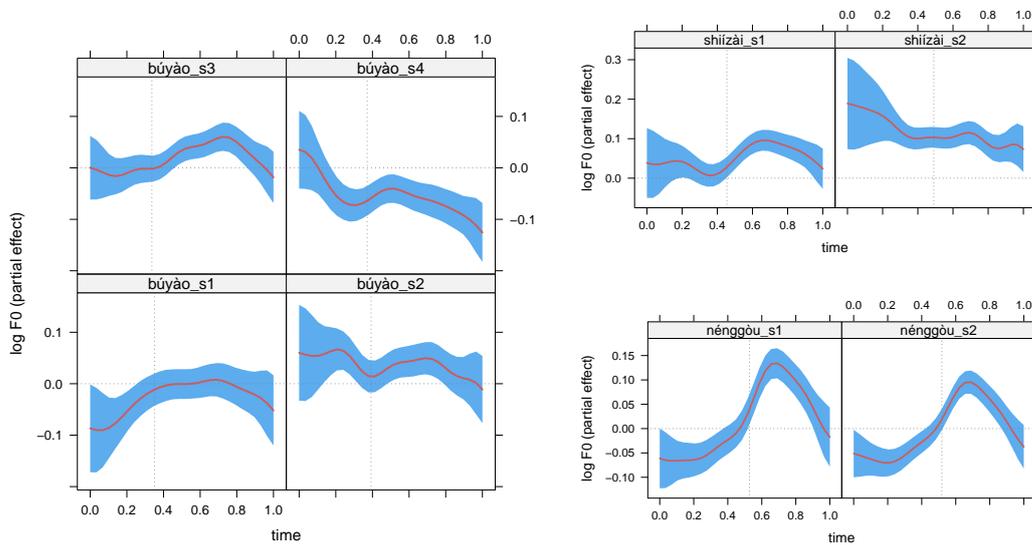


Figure 8: Sense analysis. Examples of the pitch contours predicted by the general smooth for time for female speakers, combined with the partial effects of the factor smooth for sense. The left-hand panel shows the fitted tonal contours for different senses of the word 不要 *bùyào*, a negation marker in Mandarin. The four senses are ‘prohibition’, ‘dissuasion’, ‘unnecessity’, and ‘to wish something to not happen’. The upper right-hand panel shows the fitted tonal contours for the two senses of 實在 *shízài*, meaning ‘truly’ and ‘indeed’ respectively. The lower right-hand panel plots the fitted contours for the two senses of 能夠 *nénggòu*: ‘being capable of’ and ‘enabling’.

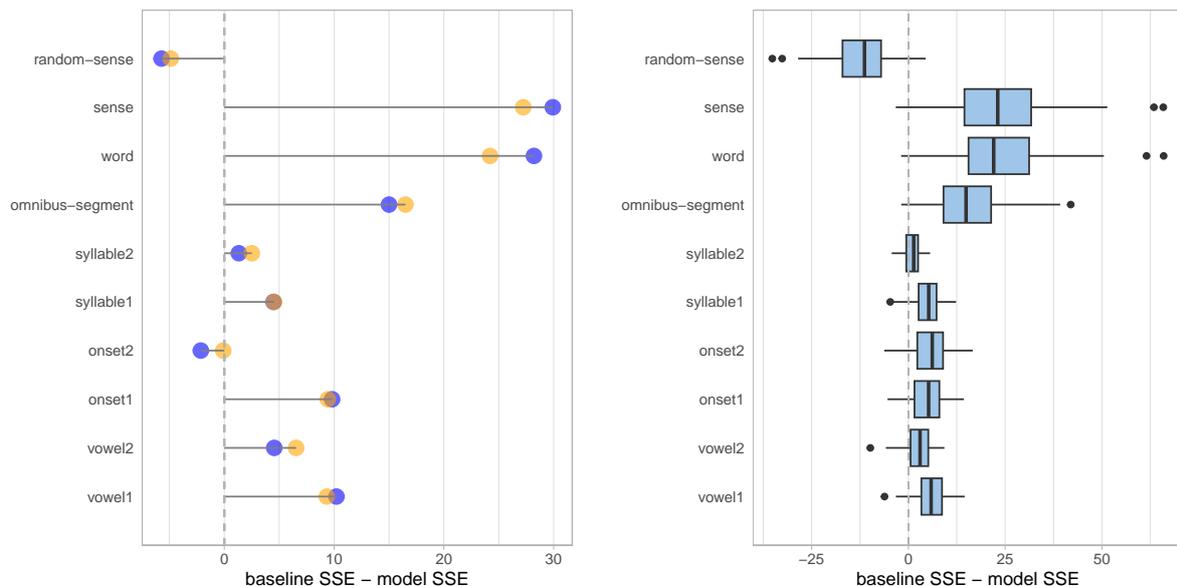


Figure 9: Model accuracy under cross validation for the sense analysis. The left-hand panel presents the reduction in summed squared errors (SSE) for models using one set of test data. Positive values indicate lower SSE as compared to the SSE of the baseline model. Orange dots represent scat models, while blue dots represent Gaussian models. The right-hand panel shows the results of 100 runs of cross validation, using Gaussian models. The boxplots represent the distributions of reduction in SSE.

2.7.2 Cross validation

As shown in the left-hand panel of Figure 9, for the first random held-out dataset the model that has access to **sense** is more accurate than the model that has access to **word**, irrespective of the type of model used (scat or Gaussian). However, in 100 Gaussian cross-validation runs, it turns out that the model with **sense** is not necessarily always more accurate than the model with **word** (right-hand panel of Figure 9). Although the median of the reduction in SSE for the sense model is slightly higher than that for the word-type model, the variance of the sense model is somewhat larger.

There are two reasons for the absence of greater prediction precision for models having access to **sense** instead of **word**. Firstly, in the smaller dataset used for these models, no fewer than 35 of the 51 word types are represented by only one sense. Any prediction advantage would therefore have to be contributed by just 16 words. Secondly, for the majority of this subset of 16 words, one sense accounts for most of the tokens. For tokens with these dominant senses, prediction is possible with greater precision. However, for tokens with less frequent senses, prediction is necessarily less precise. To see this, consider Figure 10, which presents predicted pitch contours and approximate confidence intervals for senses with many tokens (upper panels) and senses with few tokens (lower panels). Confidence bands are narrower for senses with many tokens. As a consequence, prediction for held-out tokens cannot be of the same quality for senses with few tokens as compared to senses with many tokens. The overall improvement in model fit for the sense-based GAM results from the fact that the pitch contours of tokens with the dominant sense of each word can be better predicted once these tokens are separated from tokens of minority senses.

2.8 Random forest analyses

The analyses presented thus far were all conducted with GAMs. To make sure that the effects of **word** and **sense** are not method-specific, and to further explore the relative importance of these variables in predicting F0 contours, we also made use of random forests, a machine learning algorithm that is based on multiple recursive partitioning decision trees (Breiman, 2001). Since this method focuses on overall prediction accuracy, we can include all predictors in one model, without needing to worry about concurvity or model interpretability. Figure 11 shows the variable importance scores for all our predictors in two random forest analyses carried out with the **cforest** implementation in the **party** package for R (Hothorn et al., 2006), using the same dataset as the sense-based GAM.

The left-hand panel of Figure 11 shows the results of a random forest analysis conducted straightforwardly on log F0 values. Perhaps unsurprisingly, **speaker** and **gender** emerge as the most important predictors; **sense** is ranked third and, in line with our hypothesis, contributes to model accuracy more than **word** does. However, it is surprising that among all the predictors in this model normalized **time** receives a rather low

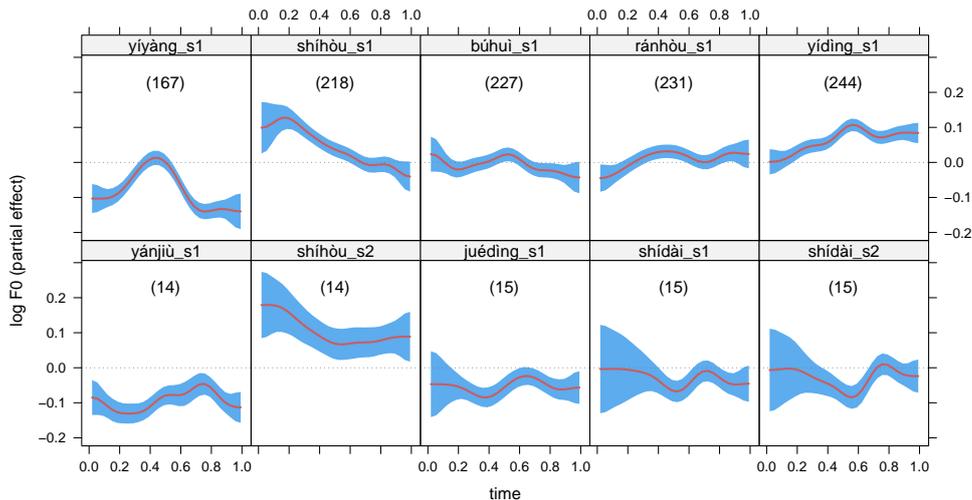


Figure 10: Sense analysis. Predicted pitch contours of the partial effects of the factor smooth for sense, for the five most frequent senses (upper row) and the least frequent senses (lower row). Numbers in the parentheses indicate the number of tokens in the dataset for the different senses. Confidence intervals are anti-conservative, due to autocorrelation in the errors.

ranking, given that the fall-rise-fall pattern is visible across the time series of many of the tokens in the dataset. Combined with the high importance of **speaker** and **gender**, this suggests that the algorithm is picking up more on the overall pitch level of tokens than on pitch variation within tokens. Furthermore, it is conceivable that individual word meanings have their own emotional valence that gives rise to higher F0 for positive words and lower F0 for negative words (Belyk and Brown, 2014).

To sidestep the confounds described in the previous paragraph, we centered the pitch contours for each speech token separately by subtracting the mean pitch of its time series from each of its pitch values. The results obtained when a random forest is requested to predict centered pitch are presented in the right-hand panel of Figure 11. Now, normalized **time** is the most important predictor, whereas **gender** has become much less important. However, **speaker** remains high in the ranking, suggesting that speaker-specificity in tonal realization is substantial. Crucially, **sense** again emerges as a more important predictor than **word**.

2.9 Summary of Section 2

The results presented in this section have provided evidence in support of our first two predictions. Firstly, variation in tonal realization cannot be reduced to the segment-related constraints on articulation previously described in the literature. Secondly, information about a word’s meaning in context improves prediction of its tonal realization in that context. To the extent that sense labels provide more fine-grained meaning distinctions than word labels do, our results suggest that meaning plays a role in shaping the realization of tonal contours in Mandarin. In other words, in addition to the relevant segmental differences previously identified in the literature, differences in meaning also contribute.

Nevertheless, as discussed in Section 2.3, sense labels impose discrete categories on semantic variation that is actually much richer, more subtle and more nebulous than can be captured by such inventories. In the computational models reported in Section 3 below, this problem did not arise since our use of the DLM enabled us to replace relatively crude sense categories with token specific semantic representations.

3 Understanding and producing item-specific F0 contours

So far we have shown that it is possible to identify meaning-specific modulations of the pitch contour for Mandarin words with the RF tonal pattern. The question therefore arises as to whether native speakers of Mandarin can in principle profit from these meaning-specific modulations. In other words, are the semantic components in words’ pitch contours sufficiently informative that they could facilitate word comprehension for the listener? A related question is whether these subtle semantic modulations are learnable for the speaker, as opposed to arising mechanically each time a word is produced, as one might expect for purely articulatory effects. As outlined in Section 1, our third and fourth predictions, repeated here for convenience, anticipate an affirmative answer to both these questions:

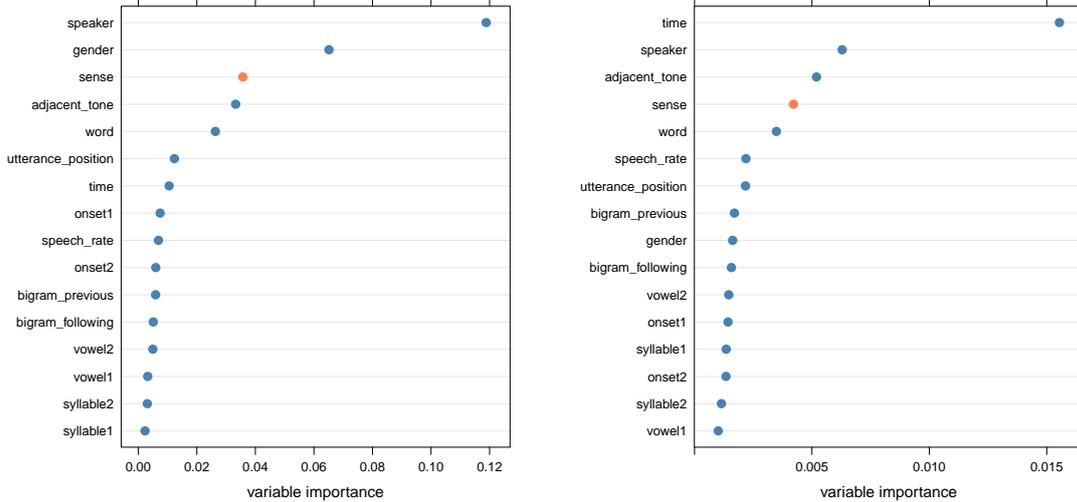


Figure 11: Variable importance scores for each predictor in the random forest analyses for log pitch (left) and centered log pitch (right).

3. Given a pitch contour, the meaning of its carrier token can be predicted above chance level, assuming the listener has previous experience of that word type.
4. Assuming they have previous experience of a given word type, a speaker can produce an appropriate pitch contour for a meaning they want to convey with that word.

In this section of the paper, we explore these predictions with computational modeling using the Discriminative Lexicon Model (DLM Baayen et al., 2019; Chuang and Baayen, 2021; Heitmeier et al., 2023c). If we can show that a simple computational model can learn to predict the meaning of a word token from its pitch contour, and that pitch contours can be predicted from intended meaning, we have a proof of concept for the functionality of meaning-specific pitch realization in human lexical processing.

As described in Section 1, the DLM focuses on the relationship between words’ forms and their meanings, and allows for fine-grained alignments between low-level features of form and low-level features of meaning. Form-meaning relationships are captured by two networks: a comprehension network that maps word form onto word meaning, and a production network that maps word meaning onto word form. Recall that in the DLM theory of the mental lexicon, forms and meanings do not have representations in memory. Form representations represent ephemeral auditory or visual input, which dynamically generates a corresponding, equally ephemeral, meaning representation. Conversely, a meaning conceptualized by a speaker at a given point in time is dynamically transformed into ephemeral representations driving articulation. In line with this theory, the DLM generates forms and meanings on the fly on a token by token basis, making it is possible to model the relationship between a given token’s specific pitch contour and that token’s context-specific meaning, and hence to account for correspondences between meaning and fine phonetic detail.¹⁶

In the DLM, both the form and the meaning of each word token are operationalized mathematically as high-dimensional numeric vectors. In order to test our predictions about the functionality of tonal modulations, the form vectors used in this study are based exclusively on the F0 contour of the relevant token, with no information about its segmental make-up. The meaning vectors that we use are context-specific, and hence also vary from token to token. Sections 3.1 and 3.2 describe how we obtained the vectors for meaning and pitch, respectively; Section 3.3 addresses the functionality of pitch in comprehension and 3.4 does the same for production.

3.1 Representing meaning: contextualized embeddings

Embeddings are widely-used numeric representations of words’ meanings, developed from the distributional semantic insight that words with similar meanings tend to occur in similar contexts (Firth, 1968; Harris, 1954; Landauer and Dumais, 1997; Salton et al., 1975). Embeddings represent word meanings as real-valued high-dimensional vectors in a semantic space (Schütze, 1992). First generation word embeddings

¹⁶This is not possible in models of speech production and comprehension that rely on stored abstract representations to mediate between form and meaning, where every token of a given word type is assumed to be associated with the same stored representations (e.g., Cutler and Clifton Jr., 1999; Levelt et al., 1999).

are static, type-level representations that model the meaning of a word type as a fixed point in semantic space, regardless of its usage in a given context. These representations therefore have difficulty distinguishing between multiple senses of a word (Pilehvar and Camacho-Collados, 2020), and although various methods have been proposed to incorporate sense or context information into type-level embeddings (see e.g., Huang et al., 2012; Iacobacci et al., 2015; Neelakantan et al., 2014; Reisinger and Mooney, 2010), most of these methods involve the use of sense-annotated corpora, which as far as we know are not available off the shelf for Mandarin and, in any case, have the disadvantage of discretizing more complex semantic variability. An alternative is to use contextualized embeddings (CEs). In contrast to static, type-level embeddings, which are based on word co-occurrences irrespective of order, CEs take into account the sequence of words in the immediate context of a target word. CEs therefore encode word meanings at the token level, and different tokens of the same word type will have different but similar context-specific embeddings (see Appendix B for an informal introduction).

To address the issue of words having context-specific meanings, this study used CEs produced for the tokens of our data by a pre-trained unidirectional language model based on the GPT-2 architecture. The model, developed by CKIP, Academia Sinica, Taiwan¹⁷, was trained on a 4.3 billion character dataset written with traditional Chinese characters. The model has 102 million parameters and encodes each character as a 768-dimensional vector. We presented the target words and their preceding contexts (consisting of all the words that occur before the target in the current utterance, as well as all the words in the immediately preceding utterance) to the GPT-2 model, and obtained two embeddings from the model, one for each character. We then averaged the two embeddings, so that every token in our dataset received a 768-dimensional vector representing its context-specific meaning.

To visualize the semantic space of the CEs, we reduced the 768-dimensional semantic space to two dimensions using tSNE (Van der Maaten and Hinton, 2008). Figure 12 shows the resulting reduced 2D plane, using convex hulls to highlight that tokens of different word types typically fall in distinct regions while tokens of the same word type form clear clusters. Although perhaps unsurprising, this distribution confirms that, despite polysemy, the word types in our data do capture a general approximation of what each token means. It is also reassuring to see that, like static embeddings, the CEs can capture inter-word semantic relations. For instance, there is a cluster of school-related words in the lower left: 學校 *xuéxiào* ‘school’, 研究 *yánjiū* ‘research’, 學到 *xué dào* ‘learn+resultative’ (see Vulić et al., 2020, for similar results).

3.2 Representing form: pitch vectors

For the DLM to implement mappings between form and meaning, every form vector input to the model has to have the same number of dimensions as all the others. However, because we took F0 measurements every 15 ms and the tokens in our data vary in duration, our tokens also vary in the number of measurement points. This means that the raw measurements cannot be used to create the form vectors. The raw pitch contours also have the problem that there are gaps due to voicelessness (cf. Appendix A, Figure A4). To overcome these two problems, we used the GAMs described in Section 2.6 to obtain smoothed pitch contours from which we could extract a standard number of measurements.¹⁸ We generated two predicted pitch contours for each token, one using predictions from the word GAM, and the other using predictions from the corresponding omnibus-segment GAM. Each of these predicted contours was then used to generate F0 predictions at 50 equally spaced time points ranging between 0 and 1 for every individual token.

Both the word GAM and the omnibus-segment GAM include all the speaker-related and context-related control variables described in Section 2.3. The only difference is that the former additionally includes **word** as the sole lexical predictor, while the latter includes six predictors specifying words’ segmental properties. Examples of GAM-generated contours (from both the word and the omnibus-segment GAMs), together with the raw F0 values, are presented in Figure 13. As can be seen, the GAM-generated contours, though generally smoothing out the undulations in raw F0s, still largely capture the overall contour shape. Moreover, since the two GAMs provide similar but not identical predicted contours, it is possible to compare their performance in the DLM. If pitch contours generated from the word GAM provide superior fits to the respective semantic vectors compared to those generated from the omnibus-segment GAM, this will provide further evidence that the **word** variable is indeed contributing some meaning-related information.

A speaker’s gender and individual characteristics such as vocal tract anatomy, idiolect, and emotional state at the time of speaking, all have strong effects both on their baseline pitch and on pitch range. Similarly, in both our word GAM and in the corresponding omnibus-segment GAM (Section 2.6), the intercepts are largely dependent on the speaker’s gender and individual identity, and differences in amplitude are largely dependent on speech rate, which we take to reflect both the speakers’ idiolect and their emotional state at the time of speaking, amongst other things. On the semantic side, in normal spoken interaction between humans,

¹⁷We use `ckiplab/gpt2-base-chinese`, which is available on <https://github.com/ckiplab/ckip-transformers>.

¹⁸The smoothing for intervocalic voiceless segments is motivated also by the consideration that the control parameters for pitch realization vary smoothly over time in the voiceless intervals.

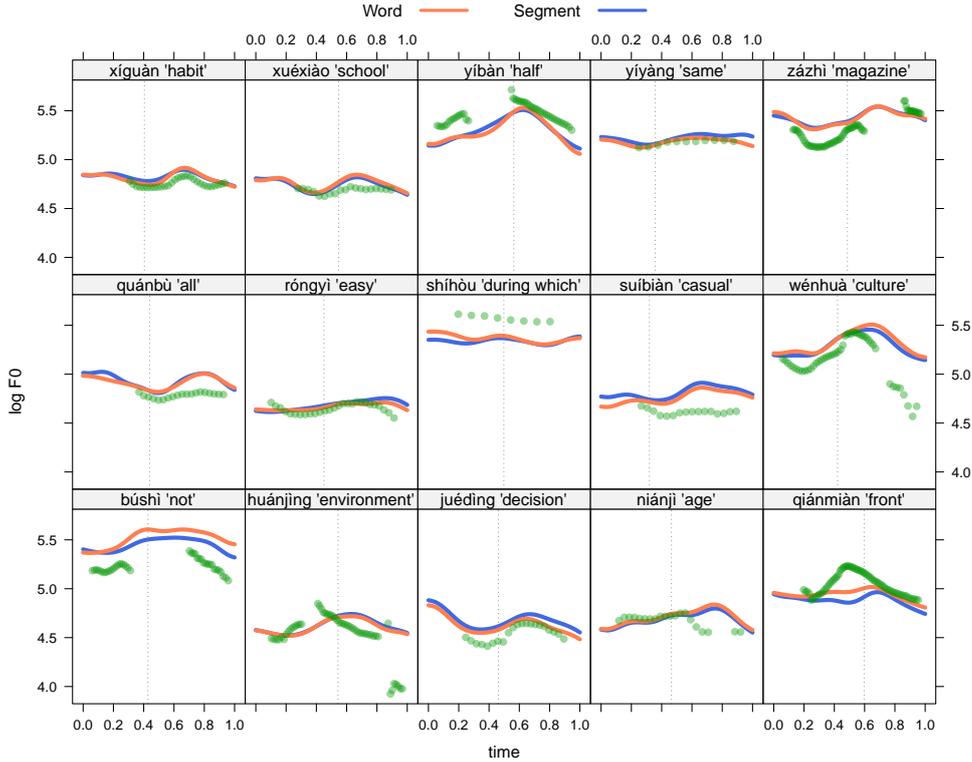


Figure 13: One token randomly selected for a selection of words. The green dots plot the observed pitch contour (raw data), and pitch vectors obtained from the word-type and the omnibus-segment models are represented by the orange and blue curves respectively. The vertical dotted lines indicate syllable boundaries.

a speaker’s identity and emotional state not only contribute to the pitch contours they produce, but are also conceptually available to their interlocutors. In contrast, the CEs used as semantic representations in our DLM modeling (Section 3.1) are based entirely on written text and therefore encode much less information about the speaker. To control for this discrepancy, we centered and scaled the predicted F0 values by token; that is, for each token in our data, and for each GAM, we calculated the mean and range of the 50 predicted F0 values, subtracted the mean from each predicted value, and divided the result by the range. In practice, this means that every token contributed equally to the model fit, irrespective of its baseline pitch or amplitude; without scaling, tokens with a greater amplitude would be taken into account more than those with a lower amplitude. A consequence of the way we centered and scaled the pitch vectors is that our DLM production models generate predictions for the geometric shapes of the contours, but not for absolute pitch or amplitude.¹⁹

3.3 Modeling comprehension

3.3.1 Method

We used two different methods to map our pitch vectors onto our semantic vectors in a comprehension network. The first method involves a straightforward linear mapping using the Linear Discriminative Learning (**LDL**) engine of the DLM. This is equivalent to the standard linear mappings used in statistics for multivariate multiple regression (see, e.g., Gahl and Baayen, 2024; Heitmeier et al., 2021, 2024, for introductions). The second method (henceforth **ResLDL**) complements the linear mapping with an additional deep mapping, making it possible to accommodate nonlinear relations while keeping the model relatively interpretable. ResLDL augments an LDL mapping with a nonlinear deep network, which is given the task of capturing any systematicities that are left unexplained in the residuals of the linear network (hence the name ResLDL). Using both these methods, and comparing the results, allowed us to shed light on the complexity of the relationship between our pitch vectors and our semantic vectors. A more detailed introduction to LDL and ResLDL is provided in Appendix C.

We split our data into a training set (80%), a validation set (10%), and a test set (10%) in such a way that every word type was represented in all three sets of data and the number of tokens per word was

¹⁹Geometric shape can be defined as ‘all the geometrical information that remains when location, scale, and rotational effects are filtered out from an object’ (Kendall, 1977).

proportional in all three sets. Both the LDL and the ResLDL mappings were trained on the training data and tested on the test data. In accordance with standard machine learning practice, the validation set was used to fine-tune the hyperparameters in the ResLDL model, before testing. This was not necessary for the LDL model since there are no hyperparameters in LDL. To ensure that our results were not specific to a particular data split, we repeated the entire modeling procedure 30 times, and therefore obtained thirty accuracy scores for each combination of pitch type (omnibus-segment or word F0 smooths) and network (LDL or ResLDL).

We evaluated the accuracy of model predictions as follows. For each pitch vector in the test set, we obtained a corresponding predicted semantic vector and identified its closest neighbour amongst the actual CEs of the tokens in our data. If this nearest neighbour belonged to any token of the same word type as the target token, the predicted semantic vector was assessed as correct, and otherwise as incorrect. This measure of success was chosen for both computational and conceptual reasons, as detailed in the following two paragraphs.

Although one might expect that a predicted semantic vector would ideally be closest in semantic space to the CE of the held-out token in question, this is computationally unrealistic. The CEs in our models are conditioned on the preceding context of a given token, and are uninformed about the following context. The pitch contours, in contrast, are shaped in part by the tone on the following word, and the probability of the word given the next word. Thus, there is information in the pitch contours that is absent in the CEs, making it computationally infeasible to predict token specific vectors. Furthermore, both the pitch contours and the CEs have measurement error. Similar to the way that a linear regression line predicts the mean value of a dependent variable for a given value of an independent variable, but not the individual data points used to generate the line, here we can predict at the level of types, but not at the level of individual tokens.

From a cognitive perspective, it is worth noting that listeners cannot arrive at exactly the same conceptualisation as the speaker, as listeners and speakers have different experiences with the language and different life histories. For example, a listener who hears ‘Do you fancy a coffee?’ may conceptualize a cappuccino, even if the speaker was envisioning an espresso. Fortunately, provided both interlocutors arrive at similar enough meanings, communication can proceed unhindered. In addition to being computationally feasible, using same word type as a criterion for success therefore makes sense in terms of human performance levels.

3.3.2 Results

Figure 14 presents the mean comprehension accuracies for the training data (left) and the test data (right). The individual barplots show accuracies for LDL (left two bars) and ResLDL (right two bars), for pitch contours based on segment-aware GAMs (red) and pitch contours based on word-aware GAMs (blue). As mentioned above, a given prediction is considered correct when the closest neighbour of the predicted CE is of the same word type as the target. For LDL, accuracy hovers around 30% for both training and test data, whereas for ResLDL accuracy is higher, over 50% and 40% for training and test data respectively. These results are surprisingly good, given that the models are requested to predict semantic vectors on the basis of pitch information only. For comparison, across our whole dataset the theoretical probability of a pitch vector and CE belonging to the same word type by chance is approximately 0.038. Similarly, baseline accuracies obtained by evaluating on a dataset with randomly permuted word labels were 3.7% for the training set, and 3.5% for the test set. This allows us to conclude that the classification accuracies of our models are far from trivial. On the contrary, even the least successful model achieves accuracies that are a whole order of magnitude greater than would be expected by chance.

A comparison of results from the two types of pitch vectors shows that meaning prediction is consistently more accurate when the pitch contour smooths are generated using the word GAM than using the omnibus-segment GAM. This indicates that the factor smooths for **word** not only contribute to a better model fit in the GAM, but also produce predicted pitch contours that are better aligned with words meanings. In other words, tonal realizations that include information about word type have the potential to help listeners to identify words’ meanings more accurately.

Finally, we note that mappings from pitch contours to CEs have significant nonlinear components, as evidenced by the higher accuracies of the ResLDL model compared to the LDL model. We further note that a nonlinear mapping may be required because we are mapping from 50-dimensional pitch contours to 768-dimensional CEs. As it is mathematically impossible to map a lower-dimensional space into a higher-dimensional space with a linear mapping,²⁰ the greater accuracy of ResLDL is unsurprising. Nevertheless, it is remarkable that the linear mappings show very similar performance on training and test data, suggesting that there is a strong linear component to predicting meaning from tonal contours.

²⁰To see this, consider points in a cube, and their projection onto a plane in that cube. From that projection, which is two-dimensional, the original locations of the points in the three-dimensional cube cannot be reconstructed.

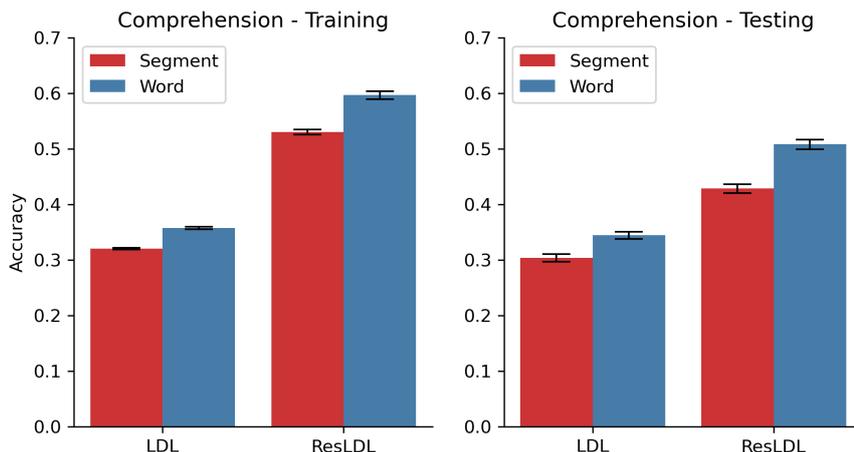


Figure 14: Mean comprehension accuracies for training data (left panel) and test data (right panel) for LDL and ResLDL mappings from omnibus-segment (red) and word (blue) pitch vectors. Mean accuracy is obtained from 30 random training and testing splits, each trained and evaluated independently. Error bars indicate double the standard error.

3.4 Modeling production

3.4.1 Method

We have seen that the pitch contours of Mandarin disyllabic words contain substantial information about word meaning. It is remarkable that a DLM comprehension model can achieve a test accuracy of over 50% when modeling with word-aware pitch contours and ResLDL. We now turn to production, addressing the question of whether a token’s pitch contour can be predicted with reasonable accuracy from its CE. If so, this would support our hypothesis that speakers can in principle learn to produce meaning-specific tonal contours.

Before going into further detail, we note that this task is considerably more difficult than the task presented to the word GAM model in Section 2. The GAM model was asked to predict pitch contours from word labels and was oblivious to variation in meaning between tokens of a given word type. In the models reported below, however, the LDL and ResLDL mappings are confronted with semantic vectors that are different from token to token. The question is whether the similarities between the CEs of tokens belonging to the same word are sufficiently consistent for the LDL and ResLDL mappings to predict appropriately similar pitch contours.

Model set-up was the same as for comprehension, except that to model production the input consisted of CEs and the output consisted of pitch vectors. We again conducted the modeling procedure 30 times. For each CE in the test set, we obtained a corresponding predicted pitch vector and identified its closest neighbour amongst the actual (GAM-generated) contours of the tokens in our data. If this nearest neighbour belonged to any token of the same word type as the target token, the predicted pitch vector was assessed as correct, and otherwise as incorrect. We complemented this overall evaluation with a qualitative analysis of the pitch contours predicted by the model for individual word types. To do this, we calculated the centroid of the CEs for all tokens of a given word, and used this centroid vector to generate a predicted pitch contour from the production network with LDL mappings and word-based pitch vectors. For each of the words presented in Figure 5, we then assessed the quality of this LDL-predicted contour by visually comparing it with the contour produced by averaging the actual (GAM-generated) pitch vectors used to train the model, for all tokens of the word in question.

3.4.2 Results

Mean production accuracies (over 30 repetitions) for the token-based evaluation are presented in Figure 15. For training data (left), accuracies are between 30% and 40%. The accuracies for the test data are only slightly lower, ranging between 27% and 35%. The probability of a CE and pitch vector belonging to the same word type by chance is the same as for the comprehension models, namely 0.038. Permutation baselines are again 3.7% for training and 3.5% for testing. In other words, like the comprehension models, the production models have accuracies an order of magnitude greater than would be expected due to chance. However, in contrast to the comprehension results, production accuracies are remarkably similar for LDL and ResLDL. Apparently, linear mappings suffice when predicting low-dimensional pitch contours from high-dimensional CEs, and succeed in capturing the regularities in the meaning-to-form mappings. Possibly,

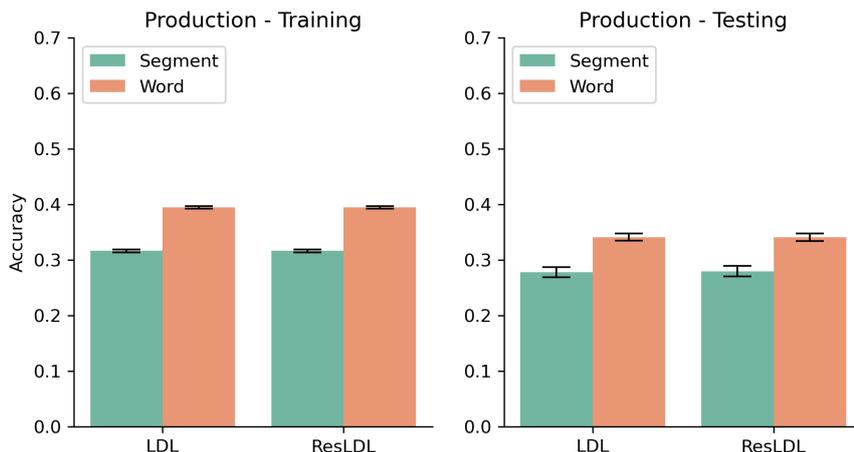


Figure 15: Mean production accuracies for training data (left panel) and test data (right panel) for LDL and ResLDL mappings from omnibus-segment (green) and word-type (orange) pitch vectors. Mean accuracy is obtained from 30 random training and testing splits, each trained and evaluated independently. Error bars indicate double the standard error.

predicting pitch from semantics is a cognitively more natural task than predicting semantics from just pitch on its own, and hence requires less powerful mappings. Finally, as for comprehension mappings, predicting pitch contours from CEs is more successful when pitch contours are generated with word-based GAMs, compared to segment-based GAMs.

The results of the qualitative analysis are shown in Figure 16. The LDL-predicted contours are shown in orange, and the by-type averages of the contours used to train the model are shown in gray. A comparison of these two contours for any given word reveals remarkable similarity, indicating that the LDL production model generates high quality predictions for the shapes of the pitch contours. It is also striking that, in shape, these orange and grey contours closely resemble the contours in Figure 5, reproduced for convenience as the blue contours in Figure 16.²¹ Recall from Section 2.6.1 that this third set of contours was produced by combining the partial effect smooth for each word type with the general smooth for time for female speakers. The similarity therefore suggests that the word-specific pitch contours isolated by our word GAM can be understood as pitch contours that correspond to the centroids of word’s contextualized embeddings. From this, we draw the conclusion that there is considerable isomorphy between the space of token-specific pitch contours and the semantic space of token-specific embeddings.

In addition to by-word centroids, we also calculated the centroid of the CEs for all tokens of all types in our dataset. The pitch contour predicted from this overall centroid is very similar to the pitch contours in the right panel of Figure 1 and the first five panels of Figure 3. In other words, the centroid of the embeddings of all tokens can be interpreted as the ‘meaning’ of the unmodulated rise-fall pitch contour. The 10 tokens that are closest to this centroid belong to the word types 不過 *búguò* ‘but, however’, 然後 *ránhòu* ‘and then’, and 時候 *shíhòu* ‘during which’, which suggests that these words are the most typical carriers of the RF pitch contour in the current dataset.

4 General discussion

This study investigated variation in the F0 contours of disyllabic words with the rising-falling (RF) tonal pattern in Taiwan Mandarin. The central hypothesis of our study is that Taiwan Mandarin disyllabic word tokens have pitch contours that are in part driven by their meanings. In standard analyses of tone in Mandarin, the rising tone of the first syllable and the falling tone of the second syllable of RF words are inherited from the single-syllable constituents and are taken to be basic, underlying tones. Deviations from these tones are explained by appealing to articulatory and prosodic constraints governing how tones can be realized. Our hypothesis adds word meaning as a missing player in the articulatory arena by arguing that meaning co-determines the realization of Mandarin tones.

Our core hypothesis generates four predictions. The first prediction is that variation in tonal realization cannot be reduced to the segment-related constraints on articulation previously described in the literature. Using generalized additive models (GAMs), we were able to show that word type is indeed a more powerful predictor of tonal realization than a wide range of words’ form properties considered jointly. We not only

²¹Note that the blue contours in Figure 16 are identical to the contours in Figure 5, except that the amplitudes are different due to the way we centered and scaled the contours in Figure 16.

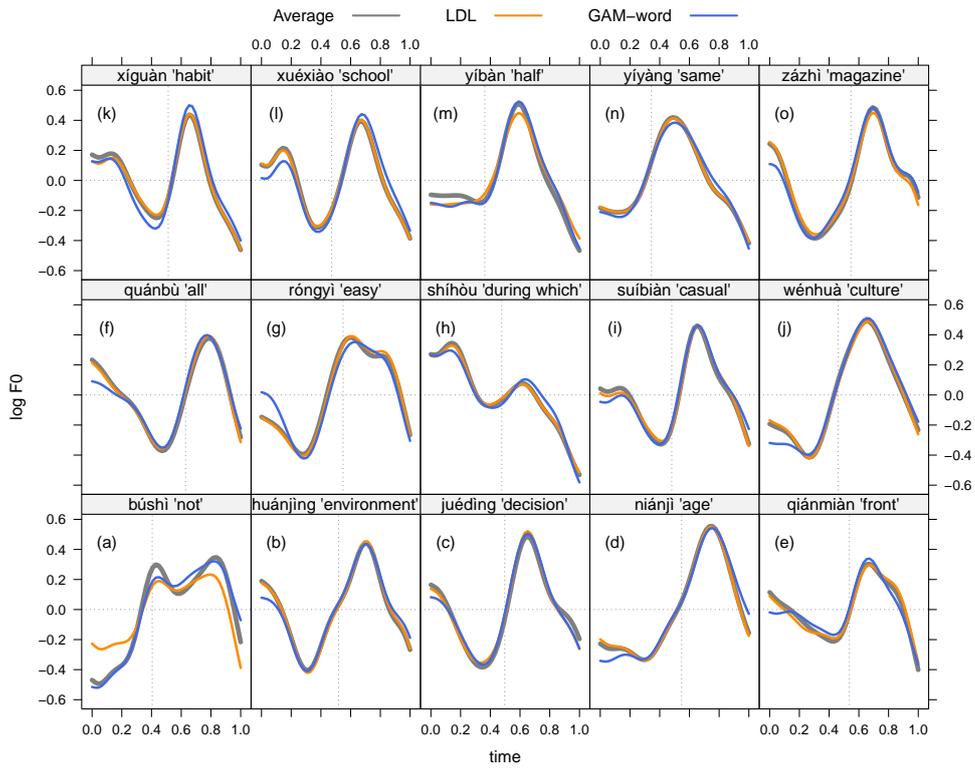


Figure 16: Examples of average pitch vectors (gray), and predictions generated by LDL (orange). The LDL contours were predicted from ‘centroid’ word meaning, obtained by averaging the CEs of all tokens of the same type. The word-specific contours predicted by the word GAM as presented in Figure 5, are reproduced here (after centering and scaling) in blue. The vertical dotted lines in the panels indicate the average (word-specific) syllable boundary.

established that GAMs with factor smooths for word type provided substantially improved model fits, but also demonstrated that word-informed GAMs provided more accurate prediction for the F0 contours of held-out data. We concluded that individual word types have specific properties — over and above their segmental makeup — that modulate the general, sine-wave shaped F0 contour characteristic of words described as having a rising-falling tone pattern. These specific properties, we conjectured, are semantic in nature.

The second prediction is that information about a word’s meaning in context will improve prediction of its tonal realization. If words with the very same segments and canonical tones, but different meanings, have distinct tonal realizations, this provides evidence for the possibility that there is a semantic component to the realization of tone in Taiwan Mandarin disyllabic words. Again using generalized additive modeling, we were able to show that adding information about meaning in context (sense) led to significant improvement in model fit. On the other hand, when it came to predicting the pitch contours of unseen tokens, the ‘word’ model performed just as well as the ‘sense’ model. Having traced this equivalence to a lack of statistical power in the sense model, we ran supplementary analyses with random forests, which again revealed greater importance for the variable including sense than for word type alone. These results provide evidence for the possibility that Mandarin disyllabic word tokens indeed have tonal realizations that are partially determined by their semantics.

The third prediction is that given a pitch contour, the meaning of its carrier token can be predicted above chance level, assuming the listener has previous experience of that word type. This prediction is important for two reasons. Firstly, for meaning-specific tonal contours to facilitate speech comprehension, it is essential that pitch contours provide the listener with information that can actually be used to reduce uncertainty about a word token’s meaning. Secondly, this prediction is also motivated by a higher-level observation. Over the centuries, the phonotactics of Mandarin have been progressively simplified, resulting in a small inventory of a mere 400 different syllables. When tones are taken into account, the number of syllables increases to 2000, which is still considerably smaller than the number of syllables in use in e.g. English, which according to a perusal of the Celex database (Baayen et al., 1995) is around 14,400. With such a restricted inventory, single-syllable Mandarin words are inevitably characterized by extensive homophony. From a functional perspective, the presence of meaning-specific pitch modulations may therefore compensate for the lack of semantic discriminability afforded by segmental makeup and syllable structure (see, e.g., Sampson, 2015, 2019, for a discussion of theoretical implications).

In order to establish empirically that token-specific tone contours indeed have the potential to help listeners zoom in on a token’s meaning, we used a specific computational modeling framework, namely the Discriminative Lexicon Model (DLM). Given the difficult task of predicting words’ high-dimensional semantic embeddings from low-dimensional pitch contours, the DLM model that we implemented performed on held-out data with an accuracy of over 50%, compared with a random baseline of 3.5%. The tonal contours of word tokens turned out to be far more revealing about their meanings than anything we thought might be possible when we started this investigation. We take these results as evidence that human listeners are presented with highly informative pitch contours, which leads us to expect that native Mandarin listeners make use of this information to optimize comprehension.

The fourth and final prediction that follows from our central hypothesis is that, assuming they have previous experience of a given word type, a speaker can produce an appropriate pitch contour for a meaning they want to convey with that word. We again tested this prediction using the DLM. A network trained on matching context-specific semantic embeddings to token-specific pitch contours performed far above a random baseline, with accuracies ranging from 30% to 40% on training data, and 25% to 30% on testing data. Given that the computational models were forced to predict pitch across tokens produced by many different speakers, without any information about the segmental make-up or syllable structure of the words, this is a remarkable result that provides strong support for actual speakers in principle being able to learn to produce meaning-specific pitch contours. At this point, however, a word of caution is appropriate. Our DLM models were given the task of predicting the geometric shape of pitch contours, and did not address token- and word-specific differences in pitch height and amplitude. The modeling of the full pitch contours, including height and amplitude, is left for future investigation.

Although the four predictions that follow from our central hypothesis are empirically well-supported, this of course does not necessarily imply that our hypothesis is correct. In general, given an implication from a proposition p (a hypothesis) to an empirical prediction q , $p \Rightarrow q$, observing q does not imply p : $q \not\Rightarrow p$, as there might be conditions other than p that also give rise to q . Thus, we cannot rule out that the importance of meaning in our models might actually be due to factors that we did not take into account in our analyses. For instance, the effects of prosody, pragmatics, syntax, and emotion could, in principle, conspire to yield effects that would seem to imply token-specific semantic effects. Measures of surprisal and informativity other than the forward and backward probabilities that we included as control variables may also be informative (Tang and Shaw, 2021). It could also be argued that in the present study, which is based on spontaneous conversational speech, the consequences of contraction and reduction (Ernestus, 2000; Johnson, 2004; Tseng, 2005) are not controlled for. And indeed, we agree, all these factors are worth

further investigation. However, it seems unlikely to us that any of these factors will turn out to explain away completely the effect of meaning. Our analyses are based on multiple tokens of each word type, that vary with respect to their syntactic position, their pragmatic function, the amount of segmental reduction, and their emotional valence. The pitch modulations estimated with the help of the GAM models are statistical generalizations across all this variation that is present in our data. It seems highly unlikely to us that the factors we were not able to control for will be distributed across our tokens in such an unbalanced way that they would be able to explain away our semantic effects. Furthermore, the simple fact that it is possible to predict embeddings from F0 contours, and F0 contours from embeddings, with an accuracy far beyond chance levels, bears witness to a non-trivial role for semantics in the tone system of (Taiwan) Mandarin Chinese.

It is an open question to what extent speakers and listeners actually can make use of the information that we have documented is present in token-specific pitch contours and token-specific embeddings. On the one hand, our computational model is trained under perfect conditions. This line of reasoning suggests that human learners are at a disadvantage compared to our computational model. On the other hand, human speakers and listeners encounter far more words than we considered in this study, and they understand or produce these words in much richer environments, including non-verbal aspects of communication such as gesture and facial expressions, than are available to our model. This alternative line of reasoning suggests that the human cognitive system is well-positioned for exploiting the isomorphisms between form and meaning to optimize comprehension and production.

A related question is whether listeners actually do make use of the distributional-statistical information that is in the speech signal of Mandarin RF words. We think they must be doing so, for two reasons. First, the pertinent information is present in the speech signal, but this information cannot be reduced to the consequences of bio-mechanical constraints on the speech production process. Second, the tone system of Mandarin Chinese is language specific, which argues against the possibility that the fine details of words' pitch contours are an inevitable consequence of their meanings. Listeners must be learning the distributional statistics of tone and meaning from the speech to which they are exposed. We hasten to note that the learning of the systematicities between pitch and semantics is in all likelihood a completely subliminal process. It is not necessary for learners to be aware of the subtle modulations of pitch contours in relation to equally subtle nuances in meaning. In our conception of the learning process, successful understanding, token by token, will drive low-level learning in the lexical networks, without conscious reflection and effort being required.²²

Our central hypothesis that pitch contours in Taiwan Mandarin disyllabic words are in part determined by semantics fits well with classical studies such as those of Bybee (2001) and Hawkins (2003).²³ These early studies are complemented with accumulating evidence that indeed the details of phonetic realization are often intimately connected to semantics. Gahl (2008) reported that English heterographic homophones tend to differ in duration, depending on their frequency, and Gahl and Baayen (2024) recently showed, using distributional semantics, that differences in homophone duration can be traced in part to the meanings of these homophones. Drager (2011) documented the ways in which the pronunciation of English 'like' varies with its meaning. Tomaschek et al. (2019a)'s study of the duration of English syllable-final /s/ likewise demonstrated the importance of semantics for phonetic detail.

Our results might be taken to suggest that word types have their own specific pitch contours stored in the mental lexicon. However, this is not what we are claiming. Within the framework of the discriminative lexicon, pitch contours are not stored. In comprehension, pitch contours contribute to predicting a word's meaning, and in production, a pitch contour is computed on the fly from the intended meaning. Because the DLM makes no attempt to derive words' forms from underlying forms, and because no abstract representations mediate comprehension or production (unlike the theories of e.g., Cutler and Clifton Jr., 1999; Levelt et al., 1999), the model provides a tool for probing the relation between form and meaning at the level of individual tokens. The facts that for comprehension the mapping from a pitch contour to its context-specific embedding is to a large extent linear, and that for production the reverse mapping is completely linear, indicate that there is considerable isomorphy between the form space of Mandarin word tokens' pitch contours and the semantic space of Mandarin words' context-specific meanings. This means that form and meaning mirror each other to a much greater extent than is often assumed, especially in frameworks that take as axiomatic that language has a 'dual articulation' (Martinet, 1965) that allocates form and meaning to two unrelated, orthogonal, components of the grammar.

Since the GPT-2 model used to generate our CEs has been pre-trained on a very large amount of data, the CEs partly capture words' statistical co-occurrence profiles. The isomorphy between pitch contours and CEs therefore suggests that the history of collocational usage affects the realization of tone, and not only the properties of the particular context (such as informativity and the properties of preceding and following

²²For token-by-token incremental lexical learning, see Heitmeier et al. (2023b), and for continuous recalibration in vision, see Marsolek (2008).

²³Note that we do not claim that pitch contours are determined exclusively by semantics; the GAM analyses that we report show ample effects of non-semantic predictors.

words) in which a token occurs. From this perspective, the present results dovetail well with conclusions reached by Seyfarth (2014) and Sóskuthy and Hay (2017), who argue that tokens of words that are often highly predictable have shorter spoken word duration even in contexts in which the factors normally driving shortening are absent (see also Tang and Shaw, 2021, for spoken word duration, maximum pitch, and maximum intensity in Mandarin).

The results that we report in the present study not only fit well with recent novel models for understanding lexical processing, but also have important consequences for teaching Mandarin as a second language. If our results generalize from the rise-fall tone pattern to other tone combinations in Mandarin disyllabic words (and ongoing research suggests that indeed this is the case), then presenting second language learners with the typical tones indicated by Pinyin transcriptions can be highly confusing and counterproductive. For instance, a learner presented with the Pinyin of 學校, *xuéxiào* ‘school’, but hearing a fall-rise-fall, and noticing the initial descending pitch, will be totally confused about what she is hearing and what (according to the Pinyin) she should be hearing. In fact, for many tokens, the feedback on the tones from the Pinyin is standing in the way of learning how Mandarin speakers actually realize tones. When error-feedback to the learner is itself error-ridden, not much progress can be expected.

To conclude, we have provided a range of observations that are consistent with the possibility that the details of how tones are realized in Taiwan Mandarin disyllabic words is partially determined by meaning in context. If our interpretation of these observations is on the right track, semantics is an important missing player in current phonetic studies of F0 modulation in Mandarin. We believe our empirical findings are sufficiently strong to open up new lines of research on the realization of pitch in tone languages. We believe that our findings will help improve pedagogical methods for teaching Mandarin to L2 learners. Last but not least, we also hope that our findings will contribute to an improved understanding of why deep learning speech processing systems are so remarkably effective and now constitute the state-of-the-art in natural language processing. Our hypothesis is that these systems can pick up systematicities between form and meaning that are not open to human introspection, but that are visible to GAMs and computational models. Crucially, we hypothesize that these systematicities are not only exploited by computational modeling algorithms, but that they are also essential, albeit subliminally, for optimizing human lexical processing in comprehension and production.

Appendix

A. Autocorrelation in pitch contours

Since F0 is a measurement of the frequency with which the vocal folds vibrate, and given the physical constraints on the muscles that control the vocal folds, it is inevitable that F0 changes only slowly with time, and that hence measurements of F0 at successive points in time are autocorrelated. For statistical evaluation, this has a far-reaching consequence: since the observations of pitch at successive points in time are not independent, unavoidably, the residuals of a standard GAM fitted to pitch contours will be also characterized by autocorrelations, violating the assumption of regression modeling that the errors should be independent and uncorrelated. As a consequence, the p-values associated with predictors in a GAM, and the confidence intervals for these predictors, will be too ‘optimistic’.

By way of example, using the current dataset, we fit a simple model predicting log-transformed F0 with speaker gender and normalized time as predictors. As shown in the left-hand panel of Figure A1, the residuals of this model are characterized by substantial autocorrelation. The auto-correlation function indicates that for this model, at lag 1, the correlation of residuals at adjacent time steps is as high as 0.95, and even at a lag of 10 timesteps, the auto-correlation is still substantial, and well above 0.6.

The **mgcv** package offers a way of capturing autocorrelations in the residuals by means of building an autoregressive process into the errors. According to this AR(1) process, the errors are generated as follows:

$$\varepsilon_t = \rho\varepsilon_{t-1} + \epsilon_i, \quad \epsilon_i \sim \mathcal{N}(0, \sigma).$$

The residual at time t is modeled as the sum of a proportion ρ of the residual at time $t - 1$ and Gaussian noise. For the present data, inclusion of an AR(1) process with $\rho = 0.95$ successfully removes most of the autocorrelation in the errors, as shown in the right-hand panel of Figure A1.

Unfortunately, this solution comes with two problems of its own. First, the autocorrelation at lag 1 varies substantially between tokens. As a consequence, a single value of ρ provides only a rough approximation that is too low for some tokens, and too high for many others. This is illustrated in Figure A2. The horizontal axis displays the distribution of by-token autocorrelation values in the errors, calculated based on the residuals of a GAM model without an AR(1) process in the errors. A large majority of tokens have high autocorrelations (median = 0.79), but there are still quite a number of tokens with much lower autocorrelations. The vertical dimension plots the autocorrelation in the residuals from a model with an AR(1) process in the errors (with $\rho = 0.95$). Although overall the autocorrelation in the errors is much lower (median = 0.29), there are still tokens with relatively high autocorrelation, and more tokens with a negative autocorrelation. Inter-token variability remains substantial. This demonstrates that the ρ parameter is set to high for some tokens, but too low for others. In other words, the problem of autocorrelated residuals cannot be properly addressed until by-token specification for the autocorrelation parameter is made possible in GAM.

Second, inclusion of an AR(1) process in the errors gives rise to residuals that do not approximate a normal distribution. For illustration, the Q-Q plots in the left and middle panels of Figure A3 display the distribution of the residuals of the models without and with an AR(1) process respectively. As can be seen, when an AR(1) process is not incorporated (left panel), residuals are already not normally distributed. They come with thick tails at both ends, indicative of being more similar to a t distribution. However, after applying the AR(1) process (right panel), the residuals exhibit a very spiky distribution, with large

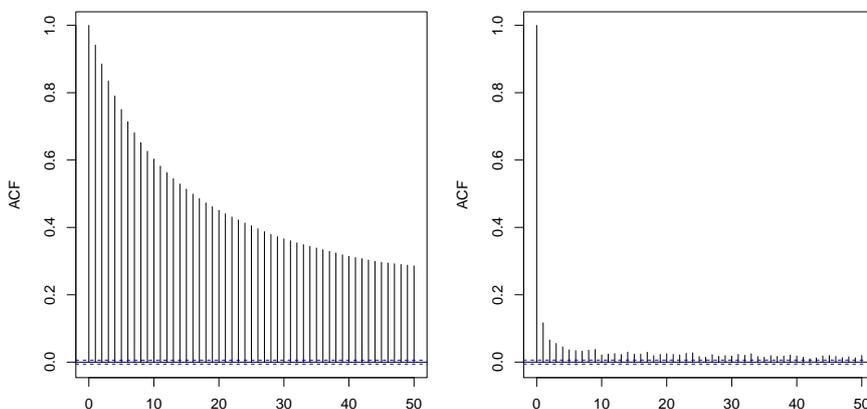


Figure A1: Autocorrelation function plots for the residuals of a GAM model without (left) and with (right) an AR(1) autoregressive process in the residuals.

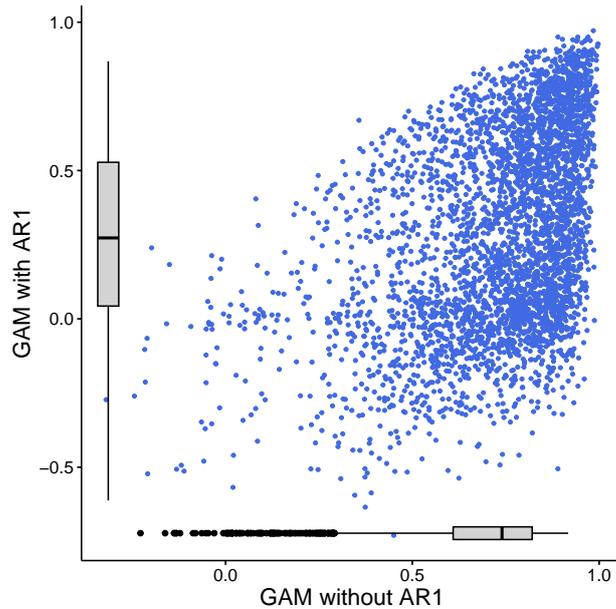


Figure A2: Scatter plot of by-token autocorrelation in the residuals from a model without AR1 incorporated (horizontal) and one with AR1 (vertical). Distributions of autocorrelation are shown by the boxplots.

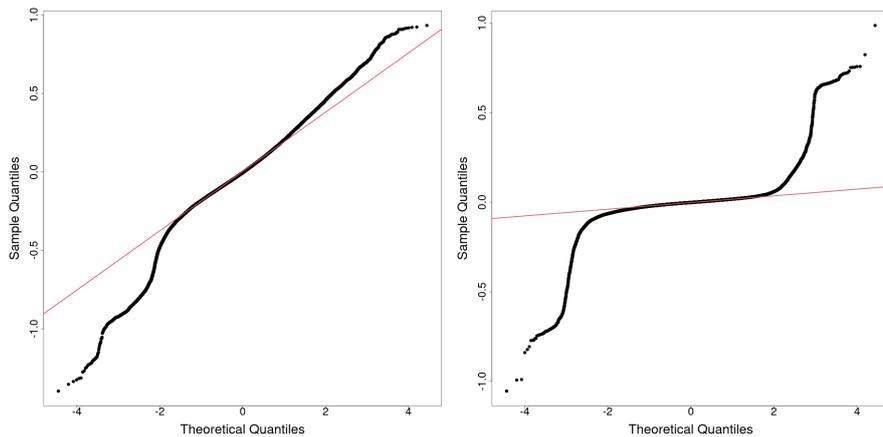


Figure A3: Q-Q plots for the residuals of a GAM model without (left) and with (right) the AR(1) process.

numbers of small errors hovering around zero, in combination with thin but very long tails. These tails resist adjustment towards normality using a scaled t-distribution (see van Rij et al., 2019, for a study where the combination of an AR(1) process and a scaled t-distribution was successful).

The large errors in the long tails of the error distribution can be traced back in part to tokens with discontinuous contours. For example, compare panel (a) of Figure A4 with panels (b)-(d). While the former shows a continuous smooth contour, the remaining three exhibit discontinuities. For panels (b)-(d), the break in the middle is due to the presence of a voiceless consonant in the onset of the second syllable. In the case of the contours in panels (c) and (d), creaky voice towards the end of the falling tone gives rise to an abrupt decrease in F0. In principle, one could start new time series after a clear discontinuity, but this raises new questions. What should counts as a ‘clear’ discontinuity? Does it make sense to start a new time series for timepoints with creaky voice? Are the pitch contours during stop closures planned smoothly but not executed, or are the planned pitch contours truly discontinuous? In the light of these considerations, we decided not to incorporate an AR(1) process in the errors.

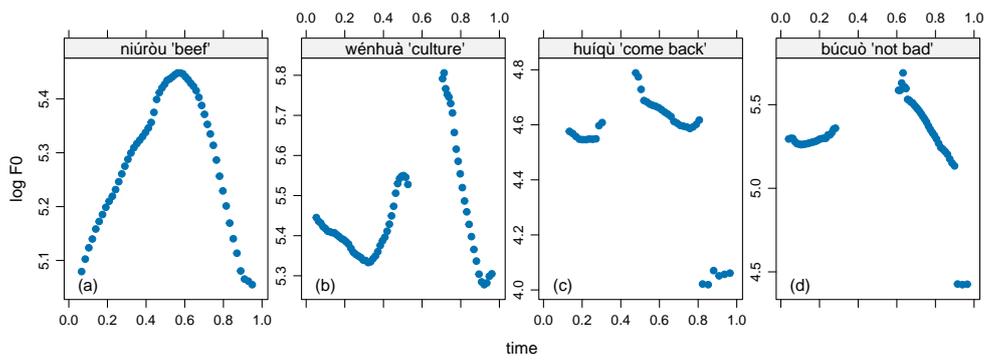


Figure A4: Panel (a) shows a complete continuous F0 contour. Panel (b)-(d) are examples of discontinuous contours.

B. Context-sensitive embeddings

Classic models of word embeddings build semantic spaces with count-based techniques. They collect the word occurrence statistics in documents or co-occurrence patterns in context windows (Landauer and Dumais, 1997; Lund and Burgess, 1996). As there are usually large numbers of documents and collocating words, the resulting embeddings are large and sparse. Therefore, dimension reduction techniques often reduce the embedding size to approximately 300.

More recent embedding models learn compact and dense representations directly from collocations. For example, Mikolov et al. (2013)’s Word2Vec model trains a single-hidden-layer neural network to predict the word given its surrounding context (continuous bag-of-words) or the surrounding context given the target word (skip-gram). The hidden representations learned by the model are the word embeddings. Different embedding models are proposed and differ in their model architecture or their prediction objectives. For instance, FastText (Bojanowski et al., 2017) uses sub-word features to construct hidden representations, and can better deal with out-of-vocabulary problems. GLoVe (Pennington et al., 2014), instead of training with CBoW or skip-gram, directly learns the word embeddings that capture global corpus statistics.

Different from the classic count-based models and the later Word2Vec or GLoVe models, contextualized embeddings (CEs) are the model representations of a deep language model (Bengio et al., 2000; Melamud et al., 2016; Raffel et al., 2020). During the training stage, the language model learns how to encode sentential contexts into CEs, given tasks such as predicting the next token or missing tokens in a sentence. Subsequently, during the inference stage, one can extract the CEs as the representations summarising the word and context in the sentence. Interestingly, although the model is not instructed to learn lexical semantics, nor are they provided with explicit sense inventories, the CEs of the same word type nevertheless reflect their context-specific or sense-level usage information (Garí Soler and Apidianaki, 2021; Pavlick, 2022; Peters et al., 2018).

Two kinds of CE-producing models can be distinguished in terms of whether the model has access to the contexts after the predicting words. In a unidirectional setting, the model is trained with a causal mask and language modeling objective, in which it learns to predict the next word only from the preceding context (e.g. GPT, Radford et al., 2018). In contrast, a bi-directional model learns through a masked language modeling task, predicting a masked word from its surrounding text (e.g. BERT, Devlin et al., 2019). Compared to bidirectional models, the embeddings of unidirectional models appear to be more aligned with human cognition in language tasks (Goldstein et al., 2022; Schrimpf et al., 2021).

C. LDL and ResLDL

Linear Discriminative Learning (LDL) is one of the computational engines of the Discriminative Lexicon Model (DLM, Baayen et al., 2019; Heitmeier et al., 2023c). It implements linear mappings between word form and meaning. Mathematically, it is equivalent to multivariate multiple regression in statistics. These linear mappings can be estimated straightforwardly with matrix algebra.

By way of example, let \mathbf{C} and \mathbf{S} denote word form and word meaning matrix respectively. The rows of the form matrix \mathbf{C} contain numeric representations of words forms, and the rows of the semantic matrix \mathbf{S} contain the corresponding semantic representations. The comprehension mapping \mathbf{F} can be obtained by solving

$$\mathbf{CF} = \mathbf{S}.$$

In the same vein, the production mapping \mathbf{G} can be obtained by solving

$$\mathbf{SG} = \mathbf{C}.$$

There are different methods of obtaining mappings \mathbf{F} and \mathbf{G} , which are documented in detail in Heitmeier et al. (2021) and Heitmeier et al. (2023a). We made use of the normal equations of regression, solving the inverse using Cholesky decomposition.

The model architecture of ResLDL is shown in Figure A5. The model consists of two routes. The route on the right is a linear mapping that implements LDL. On the left is a nonlinear network. The parallel routes of the linear and nonlinear components are similar to the residual connection model, which was first proposed for computer vision (He et al., 2016) but is also widely used in later models of natural language processing (Vaswani et al., 2017). The ResLDL takes the linear transformation matrix estimated by the least squares method as it is. This linear matrix is simply frozen, i.e., the parameters (weights) in the matrix are not updated during training. The tuneable parameters are in the nonlinear route, they are trained to capture the structure in the residuals of the linear model.

The nonlinear route is composed of three feed-forward networks (FFNs), with a sigmoid activation functions in between. The FFN itself is a linear transformation, which contains a set of tunable weights used in matrix multiplication and tunable bias terms, which act like intercepts in linear regression. The nonlinearities are introduced with the sigmoid function which maps input values into the $[0, 1]$ interval. This nonlinear function is crucial because, without it, the successive FFN will collapse into only one linear transformation.

The tuneable parameters of the ResLDL model are initialized randomly, and subsequently optimized during training to minimize the difference between the predictions and the target vectors. At first, the model does a forward pass. For example, in the case of comprehension, the input is a 50-dimensional pitch contour vector, and the model is trying to predict the 768-dimensional CE vector. The model takes a 50×768 LDL comprehension matrix for its linear route and produces a linear prediction. In addition to this linear prediction, the model also takes the pitch contour vector and inputs it into its nonlinear route. This vector is transformed by the first FFN, followed by a sigmoid function, into a latent space, which is set to 200 in this model, resulting in a 200-dimensional latent vector. This latent vector is again transformed by the following FFNs, which ultimately produce a predicted residual vector. The final output vector is the sum of the linear prediction and the predicted residual vector.

The forward pass is followed by a backward pass, using the backpropagation algorithm to update the parameters. First, the output vector is compared to the real vector, and the difference, measured by the mean squared error, is computed. Secondly, the algorithm derives the parameters' gradients from the computed difference, i.e., the directions with which the parameters should be updated to bring the predictions closer to the real vectors. Finally, the algorithm updates each parameter by subtracting the gradients multiplied by a scalar, the learning rate. The backward pass is complete with the update of all parameters. Training then proceeds with the next forward pass, using the newly updated parameters.

The model has two hyperparameters: the learning rate, and the number of epochs, i.e., the number of times the model parameters are updated on the dataset. If the learning rate is set too low, the optimisation will be too slow, while if it is set too high, we may miss the optimal parameters. Similarly, too few epochs may result in the model being undertrained, and too many iterations may lead to model overfitting. To select optimal hyperparameters, we split the data into three parts: the training set, which contains 3,022 tokens, a validation set (378 tokens), and a test set (378 tokens). The training set is used to learn the model parameters. The validation set is used to determine optimal values for the hyperparameters. Finally, the model performance is evaluated on the test set. In this study, we explored the following set of hyperparameters: two different numbers of epochs (100, 200) and six learning rates (5e-2, 1e-2, 5e-3, 1e-3, 1e-4, 1e-5). After the hyperparameter selection, the number of epochs was set to 200, and the learning rate was set to 1e-5 for the production task and 1e-2 for the comprehension task. We used the AdamW optimizer with β_1 set at 0.9 and β_2 set at 0.999, and weight decay set at 0.01. There are 200K tunable parameters in ResLDL. For each model, the training took 20 seconds on a 12th-generation Intel i7 CPU.

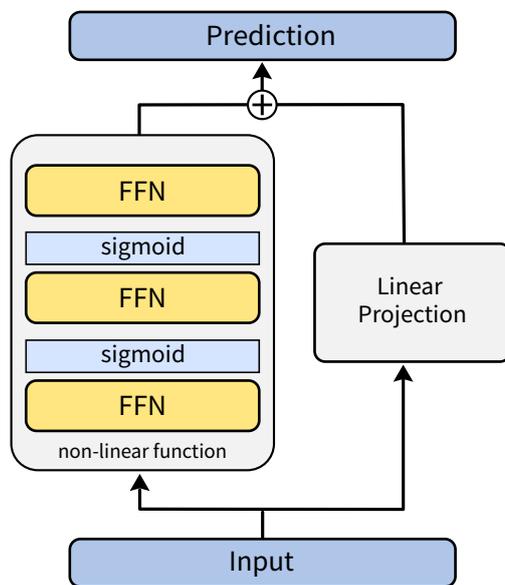


Figure A5: Model architecture of residual LDL (ResLDL). The model has two routes transforming the inputs to the predictions (either pitch contours, or contextualized embeddings). In addition to the linear route (right), which uses the standard linear transformation, the model has a second nonlinear route (left) to account for the relations not captured by the linear mapping. The nonlinear route is a three-layer fully connected neural network with sigmoid activation functions. The model adds the output from both routes to make its final predictions. Compared to the standard LDL, the ResLDL adds a nonlinear route to improve the predictions by capturing the regularities left in the residuals of linear projection.

References

- Baayen, R. H., Chuang, Y.-Y., Shafaei-Bajestan, E., and Blevins, J. P. (2019). The discriminative lexicon: A unified computational model for the lexicon and lexical processing in comprehension and production grounded not in (de)composition but in linear discriminative learning. *Complexity*, 2019:4895891.
- Baayen, R. H., Fasiolo, M., Wood, S., and Chuang, Y.-Y. (2022). A note on the modeling of the effects of experimental time in psycholinguistic experiments. *The Mental Lexicon*, 17(2):178–212.
- Baayen, R. H., Piepenbrock, R., and Gulikers, L. (1995). *The CELEX lexical database (CD-ROM)*. Linguistic Data Consortium, University of Pennsylvania, Philadelphia, PA.
- Baayen, R. H., Vasishth, S., Kliegl, R., and Bates, D. (2017). The cave of shadows: Addressing the human factor with generalized additive mixed models. *Journal of Memory and Language*, 94:206–234.
- Bell, A., Jurafsky, D., Fosler-Lussier, E., Girand, C., Gregory, M., and Gildea, D. (2003). Effects of disfluencies, predictability, and utterance position on word form variation in english conversation. *The Journal of the acoustical society of America*, 113(2):1001–1024.
- Belyk, M. and Brown, S. (2014). The acoustic correlates of valence depend on emotion family. *Journal of Voice*, 28(4):523–e9.
- Bengio, Y., Ducharme, R., and Vincent, P. (2000). A neural probabilistic language model. *Advances in neural information processing systems*, 13.
- Bi, Y., Chen, Y., and Schiller, N. O. (2015). The effect of word frequency and neighbourhood density on tone merge. In *ICPhS*.
- Boersma, P. and Weenink, D. (2019). Praat: doing phonetics by computer [computer program]. Version 6.0.48.
- Bojanowski, P., Grave, E., Joulin, A., and Mikolov, T. (2017). Enriching word vectors with subword information. *Transactions of the Association for Computational Linguistics*, 5:135–146.
- Breiman, L. (2001). Random forests. *Machine learning*, 45:5–32.
- Bybee, J. L. (2001). *Phonology and language use*. Cambridge University Press, Cambridge.
- Chao, Y. R. (1968). *A grammar of spoken Chinese*. Univ of California Press.
- Cheng, C. and Xu, Y. (2015). Mechanism of disyllabic tonal reduction in Taiwan Mandarin. *Language and speech*, 58(3):281–314.
- Chuang, Y.-Y. and Baayen, R. H. (2021). Discriminative learning and the lexicon: NDL and LDL. In *Oxford Research Encyclopedia of Linguistics*. Oxford University Press.
- Chuang, Y.-Y., Fon, J., Papakyritsis, I., and Baayen, R. H. (2021). Analyzing phonetic data with Generalized Additive Mixed Models. In Ball, M. J., editor, *Manual of Clinical Phonetics*, pages 108–138. Routledge, London.
- Chuang, Y.-Y., Huang, Y.-H., and Fon, J. (2007). The effect of incredulity and particle on the intonation of yes/no questions in Taiwan Mandarin. In *Proceedings of the 16th International Congress of Phonetic Sciences*, pages 1261–1264, Saarbrücken, Germany.
- Chuang, Y.-Y., Kang, M., Luo, X. F., and Baayen, R. H. (2023). Vector space morphology with linear discriminative learning. In Crepaldi, D., editor, *Linguistic morphology in the mind and brain*. Routledge.
- Cutler, A. and Clifton Jr., C. (1999). Comprehending spoken language: a blueprint of the listener. In Brown, C. and Hagoort, P., editors, *The Neurocognition of Language*, pages 123–166. Oxford University Press, Oxford.
- Devlin, J., Chang, M.-W., Lee, K., and Toutanova, K. (2019). BERT: Pre-training of deep bidirectional transformers for language understanding. 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 4171–4186, Minneapolis, Minnesota. Association for Computational Linguistics.
- Drager, K. K. (2011). Sociophonetic variation and the lemma. *Journal of Phonetics*, 39(4):694–707.

- Duanmu, S. (2007). *The phonology of standard Chinese*. OUP Oxford.
- Elman, J. L. (2009). On the meaning of words and dinosaur bones: Lexical knowledge without a lexicon. *Cognitive science*, 33(4):547–582.
- Ernestus, M. (2000). *Voice assimilation and segment reduction in casual Dutch. A corpus-based study of the phonology-phonetics interface*. LOT, Utrecht.
- Firth, J. R. (1968). *Selected papers of J R Firth, 1952-59*. Indiana University Press.
- Fon, J. (2004). A preliminary construction of Taiwan Southern Min spontaneous speech corpus. Technical Report NSC-92-2411-H-003-050-, National Science Council, Taipei, Taiwan.
- Fon, J. and Chiang, W.-Y. (1999). What does chao have to say about tones?-a case study of Taiwan Mandarin. *Journal of Chinese Linguistics*, 27(1):13–37.
- Fon, J. and Hsu, H.-J. (2007). Positional and phonotactic effects on the realization of dipping tones in Taiwan Mandarin. In Gussenhoven, C. and Riad, T., editors, *Phonology and Phonetics, Tones and Tunes: Vol. 2. Experimental Studies in Word and Sentence Prosody*, pages 239–269. Mouton de Gruyter, Berlin.
- Fu, J.-W. (1999). Chinese Tonal Variation and Social Network-A Case Study in Tantz Junior High School, Taichung, Taiwan. Master’s thesis, Providence University.
- Gahl, S. (2008). Time and thyme are not homophones: The effect of lemma frequency on word durations in spontaneous speech. *Language*, 84(3):474–496.
- Gahl, S. and Baayen, R. H. (2024). Time and thyme again: Connecting english spoken word duration to models of the mental lexicon. *Language*, page accepted.
- Gahl, S., Yao, Y., and Johnson, K. (2012). Why reduce? Phonological neighborhood density and phonetic reduction in spontaneous speech. *Journal of memory and language*, 66(4):789–806.
- Gårding, E. (1987). Speech act and tonal pattern in standard chinese: constancy and variation. *Phonetica*, 44(1):13–29.
- Gari Soler, A. and Apidianaki, M. (2021). Let’s Play Mono-Poly: BERT Can Reveal Words’ Polysemy Level and Partitionability into Senses. *Transactions of the Association for Computational Linguistics*, 9:825–844.
- Goldstein, A., Zada, Z., Buchnik, E., Schain, M., Price, A., Aubrey, B., Nastase, S. A., Feder, A., Emanuel, D., Cohen, A., Jansen, A., Gazula, H., Choe, G., Rao, A., Kim, C., Casto, C., Fanda, L., Doyle, W., Friedman, D., Dugan, P., Melloni, L., Reichart, R., Devore, S., Flinker, A., Hasenfratz, L., Levy, O., Hassidim, A., Brenner, M., Matias, Y., Norman, K. A., Devinsky, O., and Hasson, U. (2022). Shared computational principles for language processing in humans and deep language models. *Nature Neuroscience*, 25(3):369–380.
- Harris, Z. S. (1954). Distributional structure. *Word*, 10(2-3):146–162.
- Hawkins, S. (2003). Roles and representations of systematic fine phonetic detail in speech understanding. *Journal of Phonetics*, 31:373–405.
- He, K., Zhang, X., Ren, S., and Sun, J. (2016). Deep residual learning for image recognition. In *2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*. IEEE.
- Heitmeier, M., Chuang, Y.-Y., Axen, S. D., and Baayen, R. H. (2023a). Frequency effects in linear discriminative learning. *arXiv preprint arXiv:2306.11044*.
- Heitmeier, M., Chuang, Y.-Y., and Baayen, R. H. (2021). Modeling morphology with linear discriminative learning: Considerations and design choices. *Frontiers in psychology*, 12:720713.
- Heitmeier, M., Chuang, Y.-Y., and Baayen, R. H. (2023b). How trial-to-trial learning shapes mappings in the mental lexicon: Modelling lexical decision with linear discriminative learning. *Cognitive Psychology*, 146:101598.
- Heitmeier, M., Chuang, Y.-Y., and Baayen, R. H. (2023c). *The Discriminative Lexicon: Theory and implementation in the julia package JudiLing*. Manuscript, University of Tübingen, under review for Cambridge University Press.

- Heitmeier, M., Chuang, Y.-Y., and Baayen, R. H. (2024). *The Discriminative Lexicon: Theory and implementation in the julia package JudiLing*. in preparation for Cambridge University Press.
- Ho, A. T. (1976). The acoustic variation of Mandarin tones. *Phonetica*, 33(5):353–367.
- Hothorn, T., Hornik, K., and Zeileis, A. (2006). Unbiased recursive partitioning: A conditional inference framework. *Journal of Computational and Graphical Statistics*, 15(3):651–674.
- Howie, J. M. (1974). On the domain of tone in mandarin. *Phonetica*, 30(3):129–148.
- Hsieh, P.-j. (2013). Prosodic markings of semantic predictability in taiwan mandarin. In *INTERSPEECH*, pages 553–557.
- Hsieh, S.-K. and Tseng, Y.-H. (2020). Tutorial on sense-aware computing in chinese (version 0.1.6). In *Paper presented in 32nd conference on Computational Linguistics and Speech Processing (ROCLING 2020)*.
- Huang, C.-R., Hsieh, S.-K., Hong, J.-F., Chen, Y.-Z., Su, I.-L., Chen, Y.-X., and Huang, S.-W. (2010). Constructing chinese wordnet: Design principles and implementation. (in chinese). *Zhong-Guo-Yu-Wen*, 24:2:169–186.
- Huang, E., Socher, R., Manning, C., and Ng, A. (2012). Improving word representations via global context and multiple word prototypes. In *Proceedings of the 50th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 873–882, Jeju Island, Korea. Association for Computational Linguistics.
- Huang, P.-H. and Chiu, C. (2023). Production and perception of coarticulated tones: The cases of taiwan mandarin and taiwan southern min. *Available at SSRN 4637487*.
- Huang, Y.-H. (2008). Dialectal variations on the realization of high tonal targets in Taiwan Mandarin. Master’s thesis, National Taiwan University.
- Iacobacci, I., Pilehvar, M. T., and Navigli, R. (2015). SensEmbed: Learning sense embeddings for word and relational similarity. In *Proceedings of the 53rd Annual Meeting of the Association for Computational Linguistics and the 7th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*, pages 95–105, Beijing, China. Association for Computational Linguistics.
- Johnson, K. (2004). Massive reduction in conversational American English. In *Spontaneous speech: data and analysis. Proceedings of the 1st session of the 10th international symposium*, pages 29–54, Tokyo, Japan. The National International Institute for Japanese Language.
- Kendall, D. G. (1977). The diffusion of shape. *Advances in applied probability*, 9(3):428–430.
- Kilgarriff, A. (2007). Word senses. In Agirre, E. and Edmonds, P., editors, *Word Sense Disambiguation: Algorithms and Applications*, pages 29–46. Springer.
- Kösling, K., Kunter, G., Baayen, R. H., and Plag, I. (2013). Prominence in triconstituent compounds: Pitch contours and linguistic theory. *Language and speech*, 56(4):529–554.
- Ladd, R. and Silverman, K. E. (1984). Vowel intrinsic pitch in connected speech. *Phonetica*, 41(1):31–40.
- Landauer, T. and Dumais, S. (1997). A solution to Plato’s problem: The latent semantic analysis theory of acquisition, induction and representation of knowledge. *Psychological Review*, 104(2):211–240.
- Lee, O. J. (2005). *The prosody of questions in Beijing Mandarin*. The Ohio State University.
- Levelt, W. J., Roelofs, A., and Meyer, A. S. (1999). A theory of lexical access in speech production. *Behavioral and brain sciences*, 22(1):1–38.
- Liu, F. and Xu, Y. (2005). Parallel encoding of focus and interrogative meaning in mandarin intonation. *Phonetica*, 62(2-4):70–87.
- Lohmann, A. (2018). Cut (n) and cut (v) are not homophones: Lemma frequency affects the duration of noun–verb conversion pairs. *Journal of Linguistics*, 54(4):753–777.
- Lund, K. and Burgess, C. (1996). Producing high-dimensional semantic spaces from lexical co-occurrence. *Behavior research methods, instruments, & computers*, 28(2):203–208.
- Maaten, L. v. d. and Hinton, G. (2008). Visualizing data using t-sne. *Journal of machine learning research*, 9(Nov):2579–2605.

- Marsolek, C. J. (2008). What antipriming reveals about priming. *Trends in Cognitive Science*, 12(5):176–181.
- Martinet, A. (1965). *La Linguistique Synchronique: Études et Recherches*. Presses Universitaires de France, Paris.
- Melamud, O., Goldberger, J., and Dagan, I. (2016). context2vec: Learning generic context embedding with bidirectional LSTM. In *Proceedings of the 20th SIGNLL Conference on Computational Natural Language Learning*, pages 51–61, Berlin, Germany. Association for Computational Linguistics.
- Mikolov, T., Chen, K., Corrado, G., and Dean, J. (2013). Efficient estimation of word representations in vector space.
- Moore, C. B. and Jongman, A. (1997). Speaker normalization in the perception of mandarin chinese tones. *The Journal of the Acoustical Society of America*, 102(3):1864–1877.
- Neelakantan, A., Shankar, J., Passos, A., and McCallum, A. (2014). Efficient non-parametric estimation of multiple embeddings per word in vector space. In *Proceedings of the 2014 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 1059–1069, Doha, Qatar. Association for Computational Linguistics.
- Ouyang, I. C. and Kaiser, E. (2015). Prosody and information structure in a tone language: an investigation of mandarin chinese. *Language, Cognition and Neuroscience*, 30(1-2):57–72.
- Pavlick, E. (2022). Semantic structure in deep learning.
- Pennington, J., Socher, R., and Manning, C. D. (2014). Glove: Global vectors for word representation. In *Proceedings of the 2014 conference on empirical methods in natural language processing (EMNLP)*, pages 1532–1543.
- Peters, M. E., Neumann, M., Iyyer, M., Gardner, M., Clark, C., Lee, K., and Zettlemoyer, L. (2018). Deep contextualized word representations. In *Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long Papers)*, pages 2227–2237, New Orleans, Louisiana. Association for Computational Linguistics.
- Pilehvar, M. T. and Camacho-Collados, J. (2020). *Embeddings in natural language processing: Theory and advances in vector representations of meaning*. Morgan & Claypool Publishers.
- Plag, I., Homann, J., and Kunter, G. (2017). Homophony and morphology: The acoustics of word-final s in english1. *Journal of Linguistics*, 53(1):181–216.
- R Core Team (2022). *R: A Language and Environment for Statistical Computing*. R Foundation for Statistical Computing, Vienna, Austria.
- Radford, A., Narasimhan, K., Salimans, T., Sutskever, I., et al. (2018). Improving language understanding by generative pre-training.
- Raffel, C., Shazeer, N., Roberts, A., Lee, K., Narang, S., Matena, M., Zhou, Y., Li, W., and Liu, P. J. (2020). Exploring the limits of transfer learning with a unified text-to-text transformer. *The Journal of Machine Learning Research*, 21(1):5485–5551.
- Reisinger, J. and Mooney, R. J. (2010). Multi-prototype vector-space models of word meaning. In *Human Language Technologies: The 2010 Annual Conference of the North American Chapter of the Association for Computational Linguistics*, pages 109–117, Los Angeles, California. Association for Computational Linguistics.
- Saito, M., Tomaschek, F., and Baayen, R. H. (2023). Articulatory effects of frequency modulated by inflectional meanings. In Schlechtweg, M., editor, *Interfaces of Phonetics*. De Gruyter.
- Salton, G., Wong, A., and Yang, C. S. (1975). A vector space model for automatic indexing. *Commun. ACM*, 18(11):613–620.
- Sampson, G. (2015). A chinese phonological enigma. *Journal of Chinese Linguistics*, 43(2):679–691.
- Sampson, G. (2019). An unaddressed phonological contradiction. *International Journal of Chinese Linguistics*, 6(2):221–237.
- Schrimpf, M., Blank, I. A., Tuckute, G., Kauf, C., Hosseini, E. A., Kanwisher, N., Tenenbaum, J. B., and Fedorenko, E. (2021). The neural architecture of language: Integrative modeling converges on predictive processing. *Proceedings of the National Academy of Sciences*, 118(45):e2105646118.

- Schütze, H. (1992). Word space. In Hanson, S., Cowan, J., and Giles, C., editors, *Advances in Neural Information Processing Systems*, volume 5. Morgan-Kaufmann.
- Seyfarth, S. (2014). Word informativity influences acoustic duration: Effects of contextual predictability on lexical representation. *Cognition*, 133(1):140–155.
- Shen, X.-n. (1989). Interplay of the four citation tones and intonation in mandarin chinese. *Journal of Chinese Linguistics*, 17(1):61–74.
- Shen, X.-n. S. (1990a). *the prosody of Mandarin Chinese*, volume 118. Univ of California Press.
- Shen, X. S. (1990b). Tonal coarticulation in mandarin. *Journal of Phonetics*, 18(2):281–295.
- Shen, X. S. and Lin, M. (1991). A perceptual study of mandarin tones 2 and 3. *Language and speech*, 34(2):145–156.
- Shi, B. and Zhang, J. (1987). Vowel intrinsic pitch in Standard Chinese. In *Proceedings of the 11th International Congress of Phonetic Sciences*, pages 142–145.
- Shih, C. (1988). Tone and intonation in mandarin. *Working Papers, Cornell Phonetics Laboratory*, 3:83–109.
- Shih, C. (1997). Declination in mandarin. In *Intonation: Theory, Models and Applications*.
- Shih, C. and Kochanski, G. P. (2000). Chinese tone modeling with Stem-ML. In *Sixth International Conference on Spoken Language Processing*.
- Soskuthy, M. (2021). Evaluating generalised additive mixed modelling strategies for dynamic speech analysis. *Journal of Phonetics*, 84:101017.
- Sóskuthy, M. and Hay, J. (2017). Changing word usage predicts changing word durations in new zealand english. *Cognition*, 166:298–313.
- Sun, C. C., Hendrix, P., Ma, J., and Baayen, R. H. (2018). Chinese lexical database (cld). *Behavior research methods*, 50(6):2606–2629.
- Sun, Y. and Shih, C. (2021). Boundary-conditioned anticipatory tonal coarticulation in standard mandarin. *Journal of Phonetics*, 84:101018.
- Tang, K. and Shaw, J. A. (2021). Prosody leaks into the memories of words. *Cognition*, 210:104601.
- Tang, P. and Li, S. (2020). The acoustic realization of mandarin tones in fast speech. In *INTERSPEECH*, pages 1938–1941.
- Tomaschek, F., Plag, I., Ernestus, M., and Baayen, R. H. (2019a). Modeling the duration of word-final s in english with naive discriminative learning. *Journal of Linguistics*. <https://psyarxiv.com/4bmwg>, doi = 10.31234/osf.io/4bmwg.
- Tomaschek, F., Tucker, B. V., Ramscar, M., and Baayen, R. H. (2019b). How is anticipatory coarticulation of suffixes affected by lexical proficiency? *PsyArXiv*, pages 1–34.
- Tseng, C.-y. (1981). *An acoustic phonetic study on tones in Mandarin Chinese*. Brown University.
- Tseng, S.-C. (2005). Contracted syllables in mandarin: Evidence from spontaneous conversations. *Language and Linguistics*, 6(1):153–180.
- Van der Maaten, L. and Hinton, G. (2008). Visualizing data using t-sne. *Journal of machine learning research*, 9(11).
- van Rij, J., Hendriks, P., van Rijn, H., Baayen, R. H., and Wood, S. N. (2019). Analyzing the time course of pupillometric data. *Trends in hearing*, 23:2331216519832483.
- Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A. N., Kaiser, L. u., and Polosukhin, I. (2017). Attention is all you need. In Guyon, I., Luxburg, U. V., Bengio, S., Wallach, H., Fergus, R., Vishwanathan, S., and Garnett, R., editors, *Advances in Neural Information Processing Systems*, volume 30. Curran Associates, Inc.
- Vulić, I., Ponti, E. M., Litschko, R., Glavaš, G., and Korhonen, A. (2020). Probing pretrained language models for lexical semantics. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 7222–7240, Online. Association for Computational Linguistics.

- Whalen, D. H. and Levitt, A. G. (1995). The universality of intrinsic f₀ of vowels. *Journal of phonetics*, 23(3):349–366.
- Wieling, M., Tomaschek, F., Arnold, D., Tiede, M., Bröker, F., Thiele, S., Wood, S. N., and Baayen, R. H. (2016). Investigating dialectal differences using articulography. *Journal of Phonetics*, 59:122–143.
- Wood, S. (2017). *Generalized Additive Models: An Introduction with R*. Chapman and Hall/CRC, 2 edition.
- Wood, S. N., Pya, N., and Saefken, B. (2016). Smoothing parameter and model selection for general smooth models. *Journal of the American Statistical Association*, 111:1548–1575.
- Wu, E.-C. (2009). The Effect of Min Proficiency on the Realization of Tones and Foci in Taiwan Mandarin-Min Bilinguals. Master’s thesis, National Taiwan University.
- Wu, S.-J. (2003). A Sociolinguistic Study of Chinese Tonal Variation in Puli, Nantou, Taiwan. Master’s thesis, Providence University.
- Wu, Y., Adda-Decker, M., and Lamel, L. (2023). Mandarin lexical tone duration: Impact of speech style, word length, syllable position and prosodic position. *Speech Communication*, 146:45–52.
- Xu, C. X. and Xu, Y. (2003). Effects of consonant aspiration on mandarin tones. *Journal of the International Phonetic Association*, 33(2):165–181.
- Xu, Y. (1994). Production and perception of coarticulated tones. *The Journal of the Acoustical Society of America*, 95(4):2240–2253.
- Xu, Y. (1997). Contextual tonal variations in mandarin. *Journal of phonetics*, 25(1):61–83.
- Xu, Y. (1998). Consistency of tone-syllable alignment across different syllable structures and speaking rates. *Phonetica*, 55(4):179–203.
- Xu, Y. (1999). Effects of tone and focus on the formation and alignment of F₀ contours. *Journal of Phonetics*, 27(1):55–105.
- Xu, Y. (2001). Sources of tonal variations in connected speech. *Journal of Chinese Linguistics Monograph Series*, pages 1–31.
- Xu, Y. and Sun, X. (2002). Maximum speed of pitch change and how it may relate to speech. *The Journal of the Acoustical Society of America*, 111(3):1399–1413.
- Zhang, S., Ching, P., and Kong, F. (2006). Acoustic analysis of emotional speech in mandarin chinese. In *International symposium on chinese spoken language processing*, pages 57–66. Citeseer.