

# Hierarchical Structure and Word Strength Prediction of Mandarin Prosody

Greg Kochanski, Chilin Shih, Hongyan Jing

Bell Laboratories, Lucent Technologies

{gpk,cls,hjing}@research.bell-labs.com

## Abstract

We use Stem-ML to build an automatic learning system for Mandarin prosody that allows us to make quantitative measurements of prosodic strengths. Stem-ML is a phenomenological model of the muscle dynamics and planning process that controls the tension of the vocal folds. Because Stem-ML describes the interactions between nearby tones or accents, we were able to use a highly constrained model with only one accent template for each lexical tone category, and a single prosodic strength per word. The model accurately reproduces the intonation of the speaker, capturing 87% of the variance of  $f_0$ . The result reveals strong alternating metrical patterns in words, and shows that the speaker uses word strength to mark a hierarchy of boundaries.

## 1. Introduction

Intonation production has generally been considered a two-step process: an accent or tone class is predicted from available information, and then the accent is used to generate  $f_0$  as a function of time. Historically, most attention has been paid to the first, high level, step of the process. We here show that by focusing on  $f_0$  generation, one can build a model that starts with acoustic data and reaches far enough up to predict directly from linguistic concepts.

Specifically, we present a model of Mandarin Chinese intonation that makes quantitative  $f_0$  predictions, in terms of the lexical tones and the prosodic strength of each word. The model is able to accurately reproduce  $f_0$  in continuous Mandarin speech, with a 13 Hz RMS error. We fit this model to acoustic data and show that the strengths, tone shapes, and metrical patterns of words that result can be associated with linguistic concepts.

Further, we will show here that parameters trained on one corpus (with a properly designed model) will match equivalent parameters trained on another corpus, and also to linguistic expectations. We see effects correlated with the part of speech of words, and with the beginning and ending of the sentence, clause, phrase, and word levels of the linguistic hierarchy.

The automatic fitting is done by way of Stem-ML tags [1]. We parameterize a set of tags, then find the parameter values that accurately reproduce a training corpus.

## 2. Chinese Tones

Tonal languages, such as Chinese, use variations in pitch to distinguish otherwise identical syllables. Mandarin Chinese has four lexical tones with distinctive shapes: high level (tone 1), rising (2), low (3), and high falling (4). The syllable *ma* means *mother* with a high level tone but *horse* with a low tone. Thus, in a text-to-speech (TTS) system, good pitch prediction is im-

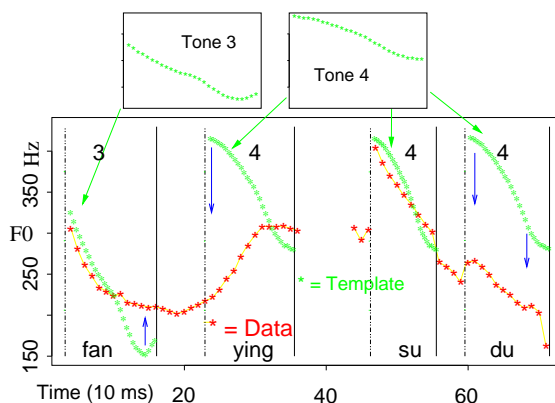


Figure 1: *Tones vs. realization.* The upper panels show shapes of tones 3 and 4 taken in a neutral environment and the lower panel shows the realization of an actual sentence containing those tones. The grey curves show the templates, and the black curve shows the  $f_0$  vs. time data.

*neutral tone*, or tone 0, which refers to special syllables with no lexical tone assignment. The pitch values of such syllables depends primarily on the tone shape of the preceding syllable.

Superficially, the modeling of Chinese tones seems straightforward. One might concatenate lexical tones to generate continuous speech. The challenge is that the realized  $f_0$  contour sometimes bears little obvious relationship to the concatenation of the tones. Figure 1 shows a Mandarin phrase *fan3 ying4 su4 du4* (“reaction time”), along with the tones from which it is constructed [2]. The last three syllables are all recognized as tone 4 by native speakers, but have drastically different  $f_0$  contours. Our model explains these changes of shape.

We explain the phenomenon displayed in Figure 1 as a natural consequence of articulatory constraints interacting with prosodic strengths. These severely distorted tone shapes occur when the shape of a weak tone is contradictory to the trajectory defined by strong neighbors. In those cases the weak tone accommodates the shapes of neighboring strong tones to maintain smooth surface  $f_0$  contours.

Our model of Chinese intonation starts with the concatenation of lexically determined tonal templates. From these, we calculate  $f_0$  at each time as a function of the nearby templates and their prosodic strengths.

Assuming that the lexical tone is known, the task of learning the Chinese prosody description given surface  $f_0$  curves

### 3. Modeling Intonation

We build our model for Mandarin on top of Stem-ML [1], because it captures several desirable properties. A positive feature of Stem-ML is that the representation is understandable, adjustable, and can be transported from one situation to another.

Unlike most engineering approaches, one can generate acceptable speech by using the templates of one speaker with parameters from another[2], where tone templates from a female speaker were used as part of a model to predict a male speaker's  $f_0$  contours. Unlike some descriptive models, we predict numerical  $f_0$  values, and so our model is subject to quantitative test, and can be extended to testing linguistic theories. Few other approaches to intonation have these properties.

Stem-ML introduces several ideas into intonation modeling:

- we assume that people plan their utterances several syllables in advance,
- we assume that people produce speech that is optimized to meet their needs,
- we apply a physically reasonable model for the dynamics of the muscles that control pitch [3], and
- we introduce a linguistically reasonable concept of a strength that is associated with each syllable.

Pre-planning in speech was first shown in terms of the control of inhaled air volume [4, 5]: people will inhale more deeply when confronted with longer phrases. This fact implies that at least a rough plan for the utterance has been constructed about 500 ms before speech begins. As another example, Figure 8 in Bellegarda *et al.* [6] shows that in an upwards pitch motion, the rate of the motion is reduced as the motion becomes longer, presumably to avoid running above the speaker's comfortable pitch range. We take this as evidence for pre-planning of  $f_0$  over a 1.5s range, at least in practiced, laboratory speech.

Next, we assume that speech is optimized for the speaker's purposes. A speaker has the opportunity to practice and optimize all the common 3-tone or perhaps 4-tone sequences, even if one assumes that each tone needs to be practiced at several distinct strength levels.

The question then arises, "optimal in what sense?" We propose that optimality be defined by a balance between the ability to communicate accurately and the effort required to communicate[1]. Specifically that the optimal pitch curve is the one that minimizes the sum of effort plus a scaled error term. Certainly, when we speak, we wish to be understood, so the speaker must consider the error rate on the speech channel to the listener. Likewise, much of what we do physically is done smoothly, with minimum muscular energy expenditure, so minimizing effort in speech is also a plausible goal.

The error term behaves like a communications error rate: it is zero if the prosody exactly matches an ideal tone template, and it increases as the prosody deviates from the template. The choice of template encodes the lexical information carried by the tones. The speaker tries to minimize the deviation, because if it becomes large, the speaker will expect the listener to misclassify the tone and possibly misinterpret the utterance.

The effort expended in speech can be approximated from knowledge about muscle dynamics [7]. Qualitatively, our effort term behaves like the physiological effort: it is zero if muscles are stationary in a neutral position, and increases as motions

related to muscle tensions. There must then be smooth and predictable connections between neighboring values of  $f_0$  because muscles cannot discontinuously change position. Most muscles cannot respond faster than 150ms, a time which is comparable to the duration of a syllable, so we expect the intonation of neighboring syllables to affect each other. Because our model derives a smooth  $f_0$  contour from muscle dynamics, our model is an extension of those of [8, 9, 10].

Effort is ultimately measured in physical units, while the communication error probability is dimensionless, so a scale factor is needed to make the two compatible for addition. This scale factor varies from syllable to syllable, and we identify it with the linguistic strength, or importance of each syllable. If a syllable's strength is large, the Stem-ML optimal pitch contour will closely approximate the tone's template, and the communication error probability will be small. In other words, a large strength indicates that the speaker is willing to expend the effort to produce precise intonation. On the other hand, if the syllable is unimportant and its strength is small, the produced pitch will be controlled by other factors: neighboring syllables and ease of production. The listener then may not be able to reliably identify the correct tone on that syllable. Presumably, the listener either can infer the tone from the surrounding context or he/she doesn't care if the listener misidentifies the tone.

We then write simple approximations to the effort and error terms, so that the model can be solved efficiently as a set of linear equations.

## 4. Experiment

### 4.1. Data Collection

The corpus was obtained from a male native Mandarin speaker reading sentences from newspaper articles, selected for broad coverage of prosodic factors. We fit two subsets (10 sentences each, 347 and 390 syllables), randomly chosen from the corpus. The speaking rate was  $4 \pm 1.4$  syllables per second, with a phrase duration of  $1.2 \pm 0.7$ s. We define phrase as speech materials separated by a pause.

Tones were identified by automatic text analysis, and checked by two native speakers. Neutral tones were manually identified. Phone, syllable, and phrase boundaries were hand-segmented, based on acoustic data.

We computed  $f_0$  with an automatic pitch tracker, then cleaned the data by hand, primarily to repair regions where the track was an octave off. If uncorrected, the octave errors would have doubled the ultimate error of the fit, and systematically distorted tone shapes.

Because word boundaries are not marked in Chinese text, different native speakers can assign word boundaries differently. Even so, the concept of a word is present, and is reflected in the prosody. We obtained word boundaries independently from three native Mandarin speakers: A, J, and S (J and S are authors). All three had generally consistent segmentation of the text into words. Pairwise comparison indicates that J and S have the highest level of agreement: J identified 395 word boundaries, S identified 370 boundaries, 99% of which were also identified by J. A identified 359 word boundaries, of which 98% agree with J's boundaries and 92% agree with S's boundaries.

Most disagreements were related to the granularity of segmentation: whether longer units were treated as single words or multiple words, and whether neutral tone syllables were at-

mented more than one way. A had the longest words, 2.04 syllables on average. J and S divided words at a finer granularity: S’s words averaged 1.98 syllables, and J’s words averaged 1.86 syllables per word. One labeler (A) consistently cliticized neutral tone syllables to the preceding word, while the other two labelers rarely did so.

We also created a random word segmentation (called “R”). The random segmentation provides a check that the metrical patterns we found are indeed significant.

## 4.2. Optimization

The Stem-ML model is built by placing tags on syllables, with adjustable parameters defining the tag shapes and positions (details below). We built several different models, focusing on models with one parameter (prosodic strength) for each word, plus a set of 36, 39, or 42 shared parameters. The models discussed here have between 210 and 246 free parameters, or an average of 0.6 parameters per syllable. The parameters that define the strength of words correlated only with a few neighbors, but the core of shared parameters are correlated with everything.

The algorithm obtains the parameters’s values by minimizing the RMS frequency difference between the data and the model. Unvoiced regions were excluded. We fit separately one the two subsets, to allow comparisons.

We used a Levenberg-Marquardt algorithm [11, 12] with numerical differentiation to find the parameters that give the best fit. The algorithm requires about 30 steps before the RMS error and parameters stabilize.

Levenberg-Marquardt, like many optimization algorithms, can become trapped in a local minimum of  $\chi^2$ , and may miss the global optimum. If we start the optimization with parameters randomly chosen from “reasonable” ranges, it will converge to what we believe to be the global minimum in about 1 in 4 tries. Consequently, we believe there are only a small number of minima. The global minimum seems to be characterized by values of *adroop*  $< 1$  (*adroop* is a Stem-ML parameter), and its  $\chi^2$  is often 10% smaller than the next best minimum. Convergence to the global minimum seems fairly reliable if an optimization is started with values of the shared parameters taken from a previous successful optimization, even if the model or data subset differ, and even if the strengths are initialized randomly.

## 4.3. Mandarin-specific Model

Our model for Mandarin is a more predictive, stronger model than bare Stem-ML [13], and is stronger even than that of [14].

The model consists of a Stem-ML *stress* tag on each syllable. We assume that each of the five lexical tone classes is described by one template. A template is defined by 5 (2 for neutral tones) pitch values, spaced across its scope. It is merely stretched (in time) and scaled (changing its pitch range) to describe all syllables which have that tone. Each tone class has a Stem-ML *type* parameter. Tone classes also have an *atype* parameter, which controls how the template scaling depends on each syllable’s strength. The pitch excursions of the template are scaled by a factor  $atype \cdot s_i^{|atype|}$  before the Stem-ML tag is generated, so that if  $|atype| > 1$ , the pitch range of the generated Stem-ML tag will change a lot for a small change in strength, while if  $|atype| < 1$ , the pitch range of the tag will be relatively independent of strength.

We give each word a *strength* parameter,  $S_w$  and derive strengths for each syllable via

where  $s_{w,i}$  is the strength of the  $i^{\text{th}}$  syllable of word  $w$ ,  $M_{L,i}$  is the metrical strength of the  $i^{\text{th}}$  position in a word of  $L$  syllables, and  $L(w)$  is the length of word  $w$ . These word strengths,  $S_w$ , are the only place in our model where linguistic information can influence the  $f_0$  contour, beyond selection of the lexical tone.

There are several parameters that are shared by all syllables. Two parameters describe the scope of templates: *ctrshift* is the offset of the template’s center from the syllable’s center, and *wscale* sets the length of the template relative to the syllable. Phrases are described by a straight-line phrase curve:

$$p(t) = P \cdot L - (D \cdot L^d) \cdot t, \quad (2)$$

where  $t$  is time,  $p(t)$  is the phrase curve, and  $L$  is the length of the phrase (in seconds). All phrase curves share three parameters:  $D$ , the declination rate;  $d$ , the dependence of the declination on the sentence length; and  $P$ , which tells how the initial height of the phrase curve depends on sentence length. To complete the model, We used Stem-ML *stepLo* tags to implement the phrase curve, and *phrase* tags were placed on phrase boundaries. Four other Stem-ML parameters control overall properties: *adroop*, *add*, *smooth*, and *base*.

We created and fit 24 different models to the data in a factorial design. We used two subsets of the corpus times the four different word segmentations (A, J, S, R) times three different parameterizations. We refer to the three parameterizations as ‘w’, ‘wA’, and ‘wAT’. These form a nested set of models with a decreasing number of parameters. In the ‘w’ parameterization, each tone class has its own *atype* and *type* parameters: we allow tone templates to scale differently as the strength increases, and we allow some tones to be defined by their shape while others are defined by their position relative to the phrase curve. In the ‘wA’ parameterization, we force all tone classes to share one *atype* parameter, so that all tone templates scale with the same function of strength. Finally, in the ‘wAT’ parameterization, we force all tones to share the *type* parameter, so all tone classes exercise the same trade-off between control of shape and control of average pitch.

# 5. Discussion

## 5.1. Results of Fit

Overall, our word-based models fit the data with a 13 Hz RMS error, approximately 1.5 semi-tones. In Figure 2, we show a typical phrase, and in Figure 3, the phrase containing the worst-fit pair of syllables in the worst model. Generally, the worst-fitting syllables tend to be the ones with the largest and fastest pitch excursions. These are conditions where Stem-ML’s approximation to muscle dynamics may break down, or where the simple approximation that we use to estimate the error between templates and the realized pitch curve may be furthest from the actual perceptual metric.

These models explain 87% of the variance of the data, and much of the rest may be explainable by phoneme-dependent segmental effects [15, 16]. Thus, essentially all the prosodic information in the  $f_0$  contour must be captured by the parameters we obtain from the fits. Of the parameters, only the word strengths have localized effects so that only they can capture localized prosodic features like emphasis, focus, and marking of sentence structure. We expect, then, that the word strengths resulting from the Stem-ML analysis are nearly a complete description of Mandarin prosody. The rest of the paper will

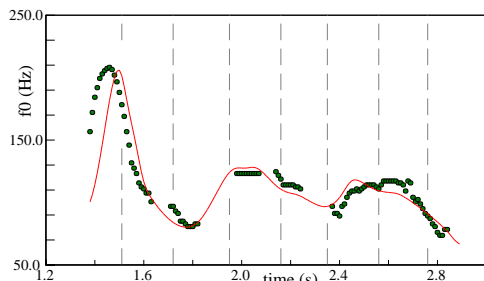


Figure 2: Typical fit (solid) vs. data (dots), for model subset1-J-A. Syllable centers are marked with vertical dashed lines.

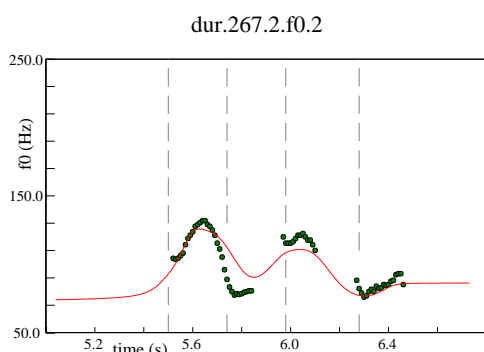


Figure 3: Phrase containing the worst-fit pair of syllables in the worst model (subset2-S-AT). Displayed as above.

We can show that the strength values that we obtain are robust against small changes in the assumptions that define the model. For example, Figure 4 shows a plot of syllable strengths obtained for the first subset with the S-wA model, plotted against strengths obtained from the J-wAT model. Despite the different word segmentations and the different sets of shared parameters the strength values are quite consistent. Comparisons between different models using the same segmentation are even closer. All the values fit on a narrow band about a smooth curve that maps the strength from one fit to the other. This mapping summarizes differences of shared parameters (most importantly *atype*) among the fits.

The strength values that are least reproducible are single syllable words, especially single syllable neutral tones.

## 5.2. Analysis of Parameters

For Stem-ML to be a model of a language, instead of just a scheme for efficiently coding  $f_0$  contours, we should be able

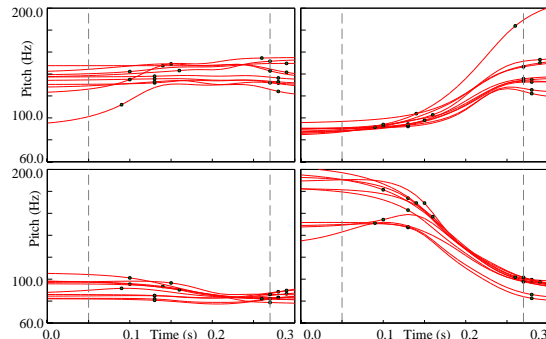
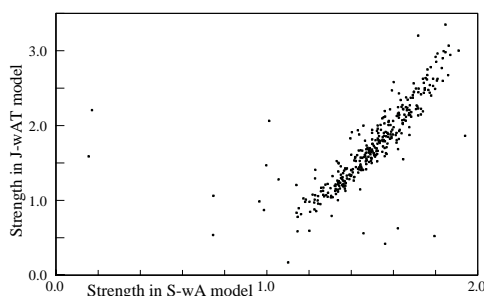


Figure 5: Modeled shapes of isolated tones. The shapes match standard descriptions, and interact to reproduce continuous speech. The two dashed vertical bars mark the syllable boundaries, and dots mark the boundaries of the tone's template in each of the 24 models. Each tone was calculated with a strength set to the median of all the strengths in the model.

to correlate the results of the fit with linguistically important features. In the following sections, we will discuss the results of the fit and see how they correlate with linguistic expectations.

Our phrase curve is Equation 2: simple linear declination. We see no evidence that the phrase curve is important, and no systematic declination. Neither  $P = -4(3) \text{ Hz}\cdot\text{s}^{-1}$  nor  $D = 0(4) \text{ Hz}\cdot\text{s}^{-1}$  is very large, and neither is substantially different from zero (error bars are shown in parentheses, and are derived from the differences between models).

In our model of Mandarin, a positive  $D$  would correspond to a systematic decrease in  $f_0$  during a phrase. This is distinguishable from a systematic decrease in strength, which would cause the magnitude of  $f_0$  swings to become smaller as the phrase progresses.

## 5.3. Analysis of Tone Shapes

First, the fitted scope of the templates is close to a syllable. The best fit templates are just 15(5)% shorter than their syllable, and their centers are offset by 18(8)% after the center of the syllable. This matches well with the intuition that tones are associated with syllables (but see [17]).

Figure 5 shows the shapes of the four main Mandarin tones in isolation, calculated for each of our 24 models. The tone shapes are consistent among different models, and across subsets. Overall, the shapes match standard descriptions of Mandarin tones. The symmetry between tones 1 and 3 and tones 2 and 4 is striking, and was in no way imposed by the analysis procedure. The four tones appear to have evolved to be nearly as different as possible.

## 5.4. Analysis of Metrical Patterns

The RMS error from these word-based models, 13 Hz, compares well with the 12 Hz RMS error we obtain from similar models[18] (with nearly twice as many parameters) that allow the strength of each syllable to vary independently, and do not impose a metrical pattern. Clearly, the metrical patterns in the words are successful at capturing much of the strength variation from syllable to syllable.

Metrical structures in words are also apparent in the fitted strengths. Figure 6 shows a tree diagram of the metrical patterns

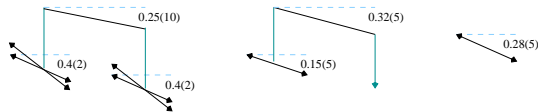


Figure 6: Metrical patterns for the J and S segmentations of 4, 3, and 2 syllable words. The words are plotted as trees, where the height of the  $i^{\text{th}}$  leaf is proportional to the metrical strength of the  $i^{\text{th}}$  syllable:  $\log(M_{L,i}) \cdot \text{atype}^{1/2}$ . Differences of  $\log(M)$  among leaves and nodes are shown numerically, with the parenthesized number showing the uncertainty in the last digit, as determined from the scatter among different models. The patterns for four syllable words have larger errors, as they are rare: they are drawn with double arrows to display the range of fitted solutions.

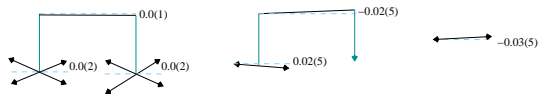


Figure 7: Metrical patterns for random word segmentation, plotted as above. As expected, the residual patterns are weak and inconsistent.

does not yield a strong metrical pattern, because there is no consistent relationship between the spoken words and the random model. Further, the R-segmentations are not as good of a fit to the data: the  $\chi^2$  for R-segmentations are 11% to 21% above the corresponding models with real (A, J, or S) segmentations. This change in  $\chi^2$  is substantial: at least an order of magnitude larger than necessary for 99% significance, even if one makes allowance for correlations among the  $f_0$  measurements.

All the real segmentations (A, J, S), show a clear strong-weak pattern for two syllable words. This means that the initial syllable’s tone is realized more precisely, and the  $f_0$  swings will tend to be larger. Although the details are strongly dependent on the circumstances, our results indicate that RMS swings on the first syllable should be 30% larger than the second syllable. While it has been generally expected that Mandarin words would show a consistent metrical pattern, previous expectations tended more to a weak-strong pattern, based primarily on evidence from duration and perceptual judgments [19].

In the A, J, and S segmentations, three-syllable words are predominantly left-branching. Because of this, we applied the same metrical pattern to all three-syllable words, and did not attempt to see if words with different internal structure had different metrical patterns. Again, we see strong-weak patterns at both levels of the metrical hierarchy, though the patterns are weaker than the two-syllable case.

All of the four-syllable words could be broken up into pairs of two-syllable words. We know this from comparison of the J and S segmentations, where the primary difference was just such a splitting and from plausibility judgments of the labelers. Consequently, we adopted the metrical tree shown in Figure 6. Expressed on that tree, we again get strong-weak patterns at both levels.

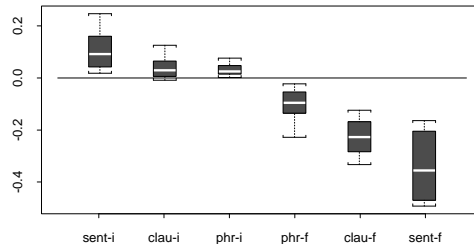
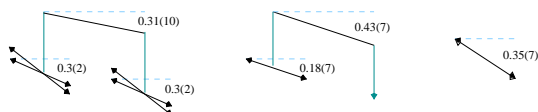


Figure 9: Correlation between strength and word positions. Each box shows the range of the data (the shaded region extends from the 25<sup>th</sup> and 75<sup>th</sup> percentiles), the median (white stripe in the box), and outlying points (brackets on the border).

In Figure 8, we show the metrical trees from the A-segmentation. While the patterns differ in detail, because of A’s tendency to attach particles to words, the pattern is similar to the J and S segmentations.

Our results are consistent with the alternating rhythmic stress patterns in Liberman and Prince [20].

### 5.5. Analysis of Word Strengths

The strengths that result from the above fitting process can be correlated with linguistically important features. We considered three features: the number of syllables in the word, the position of the word in the utterance, and the part of speech of the word, and fit the strengths with a trimmed linear regression[21] to separate out the effects of the different factors. We then ran this regression on our models, and plotted the coefficients of the factors. We found that:

(1) **Words at the beginning of a sentence, clause, or phrase have greater strengths than words at the final positions.** Figure 9 shows the regression coefficients at different positions. We define a sentence as a grammatical utterance that is marked with a period at the end, a clause as a subset of a sentence that is marked by a comma, and a phrase as a group of words that are separated by pause.

The hierarchy of linguistic units is displayed with strengths that increase with the size of the unit. Note that the zero line corresponds to the average of words that are not at a boundary, and that this line neatly divides the initial words of units from the final words of the units. These results are consistent with expectations [22].

(2) **Nouns and adverbs typically have more strength than words of other parts of speech, and particles have the lowest strengths.** Figure 10 shows the regression coefficients for different part of speech. As we can see, adverbs on average have a greater strength than words of other part of speeches. The strengths for nouns, verbs, and conjunctions are slightly weaker than that for adverbs and their strengths are close to each other. In contrast, the strength for particles (*e.g.*, neutral tones) are much weaker than that for other parts of speech.

(3) **Words with more syllables have greater strength than words with smaller number of syllables.** Figure 11 shows the regression coefficients for strengths for words of different lengths. It indicates that 3-syllable and 4-syllable words have a larger strength value than 2-syllable words, and that multi-syllable words are stronger than 1-syllable words (the 1-

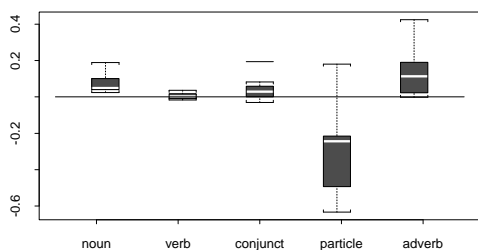


Figure 10: Correlation between part of speech and strength.

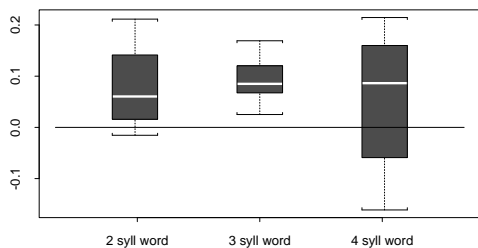


Figure 11: Correlation between strength and the number of syllables in a word.

the median absolute deviation by 17%, which (if the strength distribution were Gaussian), would correspond to Pearson's  $r = 0.31$ . We use robust estimators like a trimmed regression because the distribution of strengths has about 2% of outliers.

The correlations between strength in our Stem-ML models and the above linguistic features suggest that the strengths indeed represent the prosody importance of syllables and words. On one hand, we can use the strengths from Stem-ML models to test linguistic theories; on the other hand, we can use features such as position, part of speech, and number of syllable in word to predict the strength of a word, and thus improve prediction of  $f_0$ .

## 6. Conclusion

We have used Stem-ML to build a model of continuous Mandarin speech that connects the acoustic level to text analysis results (part-of-speech information, and word, phrase, clause, and sentence boundaries). When fit to a corpus, the model shows that prosody is used in a consistent way to mark divisions in the text: sentences, clauses, phrases, and words all start strong and end weak. Our prosodic measurements also show a useful correlation with word length and the part of speech of words.

The simplicity and compactness with which one can describe Mandarin using this representation suggests that it captures some important aspects of human behavior during speech. For more information, see <http://www.bell-labs.com/project/tts/stem.html>.

## 7. References

- [1] Greg P. Kochanski and Chilin Shih, "Stem-ml: Language independent prosody description," in *Proceedings of the*
- [2] Chilin Shih and Greg P. Kochanski, "Chinese tone modeling with stem-ml," in *ICSLP*, Beijing, China, 2000.
- [3] H. Hollien, "In search of vocal frequency control mechanisms," in *Vocal Fold Physiology: Contemporary Research and Clinical Issues*, Diane M. Bless and James H. Abbs, Eds., pp. 361–367. College-Hill Press, San Diego, CA, 1981.
- [4] Carol N. Wilder, "Chest wall preparation for phonation in female speakers," in *Vocal Fold Physiology: Contemporary Research and Clinical Issues*, Diane M. Bless and James H. Abbs, Eds., pp. 109–123. College-Hill Press, San Diego, CA, 1981.
- [5] A. L. Winkworth, P. J. Davis, R. D. Adams, and E. Ellis, "Breathing patterns during spontaneous speech," *Journal of Speech and Hearing Research*, vol. 38, no. 1, pp. 124–144, 1995.
- [6] J. Bellegarda, K. Silverman, K. Lenzo, and V. Anderson, "Statistical prosodic modeling: from corpus design to parameter estimation," *IEEE Transactions on Speech and Audio Processing*, vol. 9, no. 1, pp. 52–66, 2001.
- [7] K. N. Stevens, *Acoustic Phonetics*, The MIT Press, 1998.
- [8] S. Öhman, "Word and sentence intonation, a quantitative model," Tech. Rep., Department of Speech Communication, Royal Institute of Technology (KTH), 1967.
- [9] Hiroya Fujisaki, "Dynamic characteristics of voice fundamental frequency in speech and singing.," in *The Production of Speech*, P. F. MacNeilage, Ed., pp. 39–55. Springer-Verlag, 1983.
- [10] Yi Xu and X. J. Sun, "How fast can we really change pitch? maximum speed of pitch change revisited," in *ICSLP*, 2000.
- [11] K Levenberg, "A method for the solution of certain problems in least squares," *Quart. Applied Math.*, vol. 2, pp. 164–168, 1944.
- [12] D. Marquardt, "An algorithm for least-squares estimation of nonlinear parameters," *SIAM J. Applied Math*, vol. 11, pp. 431–441, 1963.
- [13] Greg Kochanski and Chilin Shih, "Soft templates for prosody mark-up," Tech. Rep., Bell Laboratories, Lucent Technologies, <http://www.bell-labs.com/project/tts/stem-MLdefine.pdf>, 2001.
- [14] Greg P. Kochanski and Chilin Shih, "Soft templates for prosody mark-up," *Submitted to Speech Communications*, 2001.
- [15] W. Lea, "Segmental and suprasegmental influences on fundamental frequency contours," in *Consonant Types and Tones*, L. Hyman, Ed., pp. 15–70. University of Southern California, Los Angeles, 1973.
- [16] Kim E. Silverman, *The Structure and Processing of Fundamental Frequency Contours*, Ph.D. thesis, University of Cambridge, 1987.
- [17] Yi Xu, "Pitch targets and their realization: Evidence from Mandarin Chinese," *Speech Communication*, vol. 33, pp. 319–337, 2001.
- [18] Greg Kochanski and Chilin Shih, "Automated modelling

- [19] Mao-Can Lin and Jingzhu Yan, “The stress pattern and its acoustic correlates in Beijing Mandarin,” in *Proceedings of the 10th International Congress of Phonetic Sciences*, 1983, pp. 504–514.
- [20] M Y. Liberman and A. Prince, “On stress and linguistic rhythm,” *Linguistic Inquiry*, vol. 8, pp. 249–336, 1977.
- [21] MathSoft, Inc., *Splus online documentation*, 3.3 edition, 1995, Subroutine *ltsreg()*, set to exclude the 5 most extreme data from the objective function.
- [22] Julia Hirschberg and Janet Pierrehumbert, “The intonational structuring of discourse,” in *Proceedings of the 24th Annual Meeting of the Association for Computational Linguistics*, 1986, vol. 24, pp. 136–144.