

# Discussion

## Spoken English Intelligibility Remediation with PocketSphinx Alignment and Feature Extraction Improves Substantially over the State of the Art

Yuan Gao  
17zuoye.com and  
Beijing University of Technology<sup>1</sup>

Brij Mohan Lal Srivastava  
International Institute of Information  
Technology, Hyderabad, India<sup>2</sup>

James Salsman  
17zuoye.com and  
TalkNicer.com, LLC<sup>3</sup>

*Automatic speech recognition is used to assess spoken English learner pronunciation based on the authentic intelligibility of the learners' spoken responses determined from deep neural network (DNN) model predictions of transcription correctness. Using numeric features produced by PocketSphinx alignment mode and many recognition passes searching for the substitution and deletion of each expected phoneme and insertion of unexpected phonemes in sequence, the DNN models achieve 97% agreement with the accuracy of Amazon Mechanical Turk crowdworker transcriptions, up from 75% reported by multiple independent researchers. Using such features with DNN prediction models can help computer-aided pronunciation teaching (CAPT) systems provide intelligibility remediation. We have developed and published free open source software so that others can use these techniques.*

### 1. Introduction

**Authentic intelligibility**, the ability of listeners to correctly transcribe recorded utterances, initially used for CAPT by [Kibishi and Nakagawa 2011](#) and [Kibishi, Hirabayashi, and Nakagawa 2015](#), is a better measure of pronunciation assessment for spoken language learners compared to mispronunciations identified by expert pronunciation judges or panels of experts, because such mispronunciations are associated with only 16% of intelligibility problems, according to [Loukina et al. 2015](#), who state:

We investigated . . . which words are likely to be misrecognized and which words are likely to be marked as pronunciation errors. We found that only 16% of the variability in word-level intelligibility can be explained by the presence of obvious mispronunciations. Words perceived as mispronounced remain intelligible in about half of all cases. At the same time . . . annotators were often unable to identify the word

---

<sup>1</sup> gyfreedom93@163.com (corresponding author for inquiries in Chinese)

<sup>2</sup> contactbrijmohan@gmail.com (corresponding author for inquiries in Indic languages)

<sup>3</sup> jim@talknicer.com (corresponding author for inquiries in other languages)

Submission received: August 24, 2017.

when listening to the audio but did not perceive it as mispronounced when presented with its transcription.

This substantial improvement is not yet well understood by most CAPT researchers and commercial software publishers. Currently, expert human pronunciation judges assess student performance, often with large inter-rater variability between experts scoring the same utterances. Since most formal mispronunciations do not substantially impede understanding of spoken language, automatic speech recognition CAPT systems trained to approximate the subjective assessments of judges do not perform as well as might be expected after intensive work on the issue by several hundred researchers spanning decades (Chen and Li 2016, Llisterri 2016.) While there are many commercial CAPT applications, there is no consensus among speech language pathologists about which of them, if any, work well (Dudy et al. 2017.)

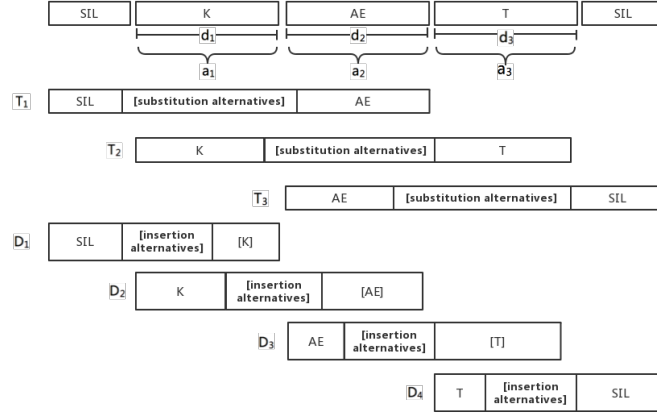
In high stakes situations, systems imitating subjective assessments of human judges have, for example, prevented native English speakers and trained English language radio announcers from immigrating to Australia (Ferrier 2017, Australian Associated Press 2017.) A more technical related problem with traditional CAPT approaches is that popular pronunciation assessment metrics, primarily **goodness of pronunciation** (GOP) as defined by e.g., Witt and Young 2000, are quotients with such vaguely specified denominators (Qian, Meng, and Soong 2015) that they tend to correlate weakly with authentic intelligibility. Earlier work suffers from similar problems.

We are offering remediation of authentic intelligibility for English CAPT to 17zuoye.com's 30 million K-6 English language students in China, and we are deploying the same technology in the Wikimedia Foundation's **Wiktionary** dictionaries along with their phonetics and pronunciation articles in **Wikipedia** to provide free CAPT assessment and remediation exercises. We are measuring which feedback choices perform the best for student proficiency outcomes, and studying the possibility of using students to provide transcriptions instead of paid **crowdworkers**.

## 2. Adapting PocketSphinx for feature extraction

We chose to use Carnegie Mellon's **PocketSphinx** free open source automatic speech recognition (Huggins-Daines et al. 2006) system's alignment routines after initially trying PocketSphinx for two-pass alignment to a fixed grammar by using the time endpoints from recognizing the phonemes of the expected utterance in sequence, using a finite state grammar with no alternative or optional components other than silence, defined using a **JSpeech Grammar Format** file. The results for the first pass were discarded, because its purpose was solely to perform **cepstral mean normalization** for adapting to the audio characteristics of the microphone, channel, and noise. We found that grammar-based alignment, which is optimized for speed instead of accuracy, resulted in less correctly predictive features than using a single pass of the alignment API functions, which are only available from the PocketSphinx C API instead of command line invocations. Using the alignment API's C functions, we were able to obtain 97% accurate prediction of authentic intelligibility, while the grammar-based alignment provided only 92% accuracy using an identical DNN classifier.

The results of the alignment are used to select audio sub-segments of the utterance to indicate substitutions of expected phonemes along with insertions and deletions, in multiple subsequent recognizer passes of each three and two adjacent phonemes at a time, respectively.

**Figure 1**

Feature extraction: The three phonemes of the word ‘cat’ are aligned, producing durations  $d_n$  and acoustic scores  $a_n$ . Then several passes of recognition to the audio aligned to groups of three ( $T_n$ ) and two ( $D_n$ ) phonemes are used to measure phoneme substitutions, and insertions and deletions, respectively.

After alignment, we run the recognizer on each sub-segment of the audio corresponding to each three aligned phonemes in sequence, and count how soon the expected phoneme occurs in the **n-best recognition results**. Then we run the recognizer on each sub-segment corresponding to each two adjacent phonemes in sequence, simultaneously counting how frequently the initial expected phoneme is omitted when searching for the insertion of all 39 phonemes and silence in between the two expected phones.

The substitution detection pass focuses on three adjacent phonemes at a time as located by the alignment routine. For the audio sub-segment of each three adjacent phonemes from the alignment, we use a grammar specifying the first and last of the three as the only options on the ends, with an alternative allowing for any one phoneme (including diphthongs) in the middle. The score, in the range  $[0, 1]$ , represents how high the expected middle phoneme ranks in the n-best results of all the possible phonemes in between the other two. We ask the recognizer for as many n-best results as possible, because sometimes a truncated grammar result (e.g., only two phonemes instead of three) result, but we often get at least 30 results from the 40 possible phonemes and silence, and sometimes get 70 results. The insertion and deletion pass operates on the audio sub-segments of two adjacent phonemes at a time, using a grammar to look for the first expected phoneme in the front as the only possibility, followed by an optional alternative of any phoneme other than the expected second phoneme counting as insertions, and then followed by the expected second phoneme specified as optional to account for deletion. Each time an insertion or deletion is returned in the n-best results before only the expected two phonemes are returned, the  $[0, 1]$  score is reduced.

We also produce each phoneme’s duration and the logarithm of its acoustic score from the alignment phase as features in our DNN classifier feature inputs. For each phoneme, we produce: (1) a duration; (2) an acoustic score from the alignment, corresponding to the numerator of the GOP score of [Witt and Young 2000](#); (3) a  $[0, 1]$  score measuring phoneme substitution, and (4) a  $[0, 1]$  score measuring insertions and

deletions. One final additional insertion and deletion measurement appears at the end of the feature vector for each word; in a multi-word phrase, that final score is shared as identical to the first insertion and deletion measurement of the next word.

We use some non-standard PocketSphinx parameters. We use a frame rate of 65 frames per second instead of 100, because learners are not likely to speak very quickly. We use a **-topn** value of 64 instead of 2. This provides more accurate recognition results at the expense of longer runtime, but our feature extraction system runs in better than real time in a single thread of a 2016 Apple MacBook Air, and on user's browsers as a **pocketsphinx.js** adaptation in **JavaScript**. We use a **-beam** parameter of  $10^{-57}$ , a **-wbeam** parameter of  $10^{-56}$ , and a **-maxhmmfp** value of  $-1$  for the same reason. We set **-fsgusefiller** to "no" so that optional pauses are not assumed between every word, allowing us to define words comprised of a single **CMUBET** phoneme without slowdown.

## 2.1 Compiling featex.c with PocketSphinx

The C source code to perform the feature extraction, **featex.c**, and instructions for compiling and using it are available under the MIT open source code license at:

<https://github.com/jsalsman/featex>

## 3. Using pocketsphinx.js in web browsers

Feature extraction can take place in web browsers' JavaScript code using the **Emscripten** system of compiling C to JavaScript, and audio recorded in web browsers supporting microphone input. During the initialization process, the browser is checked for microphone availability and the sampling frequency at which it operates. A media source stream is requested to record audio from the microphone, and connected to a recorder thread which listens or stops listening based on browser user interface events. The **pocketsphinx.js** module is initialized inside a web worker to asynchronously call the alignment and feature extraction modules.

### Web client algorithm:

**Step 1:** The user presses the 'Record' button.

**Step 2:** The recorder thread starts listening.

**Step 3:** The user presses the 'Stop' button.

**Step 4:** The recorded audio is converted and downsampled if necessary.

**Step 5:** The extracted feature vector and word sent to the intelligibility prediction service (see Section 7.)

**Step 6:** Assessment feedback is provided to the user.

The integrated code and detailed compilation instructions can be found at:

<https://github.com/brijmohan/pocketsphinx.js>

For more information and an example of an integrated web browser system, please see:

<https://github.com/brijmohan/iremedy>

For an example of how such a system might be integrated into Wiktionary, please see:

[https://brijmohan.github.io/iremedy/single\\_line.html](https://brijmohan.github.io/iremedy/single_line.html)

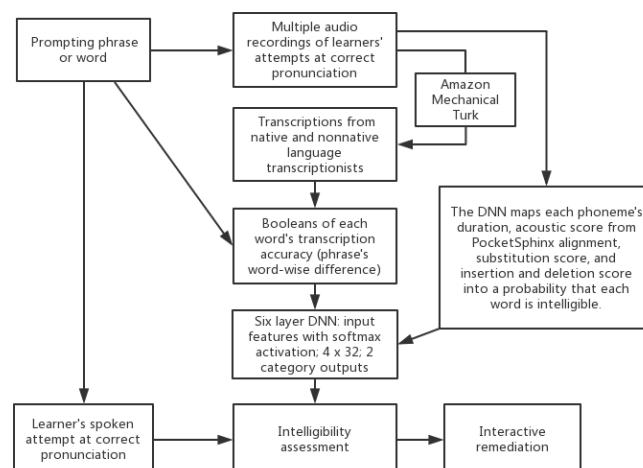
#### 4. Obtaining transcriptions of student utterances

When measuring intelligibility from transcriptions of student utterances, we count homophones obtained from the **CMUDICT** phonetic dictionary used by PocketSphinx as intelligibly transcribed.

We consistently obtained faster responses from **Amazon Mechanical Turk** when paying \$0.03 per transcript compared to \$0.15. We believe crowdworkers prefer to do low-paying tasks because they are likely to be easier and will cause fewer problems if the work is rejected. We are studying the possibility of using our English learners to provide transcriptions instead of paying crowdworkers, as *bona fide* listening comprehension and typing exercises suitable for assessments in their own right

#### 5. Predicting intelligibility

Our DNN model is defined in **Python** with **Keras** as an input layer with one unit for each of the features described in the preceding section for each of one or more words in a prompt for the learners to speak. We stack 4 more network layers of 32 units each, plus a two unit classifier output layer, using 25% **dropout** to prevent over-fitting, **softmax** activation on the input and output layers to properly handle probabilities in  $[0, 1]$ , the **ADAM** optimizer, the **glorot\_uniform** initialization algorithm, and one thousand **training epochs**. With that DNN configuration, we obtain 97% accuracy in predicting the intelligibility of 10,166 transcripts of 2,337 recordings of 82 different basic English



**Figure 2**  
Predicting intelligibility.

words. We have obtained similar results on longer phrases. Using the same features to train a **linear logistic regression** model, we only get 75% accuracy, which was reported by Kibishi and Nakagawa 2011 and Kibishi, Hirabayashi, and Nakagawa 2015, and the Educational Testing Service (Loukina et al. 2015.) Using the same features with a **support vector machine** classifier configured with a **radial basis function kernel**, we obtain only 76% accuracy predicting intelligibility from the same features.

### 5.1 Python Keras DNN specification

```
from keras.models import Sequential
from keras.layers import Dense, Dropout
from keras.utils.np_utils import to_categorical
from numpy import asarray

layers = 4; units = 32; epochs = 1000; drop = 0.25
model = {} # empty dictionary

# NOT SHOWN: set 'word' and its number of 'features'

model[word] = Sequential() # DNN
model[word].add(Dense(units, input_dim=features,
    activation='softmax',
    kernel_initializer='glorot_uniform'))
model[word].add(Dropout(drop))
for i in range(layers):
    model[word].add(Dense(units,
        kernel_initializer='glorot_uniform'))
    model[word].add(Dropout(drop))
model[word].add(Dense(2, activation='softmax',
    kernel_initializer='glorot_uniform'))
model[word].compile(optimizer='adam',
    loss='categorical_crossentropy')

# NOT SHOWN: Read feature array 'X[n][f]' and
# corresponding boolean intelligibility vectors 'y[n]'.
# n: number of transcriptions; f: number of features

y_cat = to_categorical(y) # 0: unintelligible or 1: good

model[word].fit(X, y_cat, epochs=epochs, verbose=0)
# X and y are a list of vectors and categorical booleans

# now you can get the probability of intelligibility for
# some featex.c output vector Z[f] this way:
p_i = model[word].predict(asarray(Z).reshape(1, -1))[0][1]
```

For a client-server system to predict word intelligibility from feature vectors, please see:

<https://github.com/brijmohan/proneval-service>

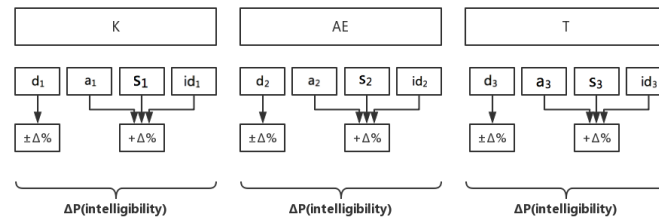
## 6. Measuring the accuracy of intelligibility assessment

When different transcripts of the same utterances result in differing intelligibility results, we measure accuracy as a fraction of the best possible result. For example, if the same utterance was transcribed correctly by three transcriptionists but incorrectly by a fourth, the maximum unadjusted accuracy achievable from predicting that utterance's intelligibility is 75%, so an unadjusted accuracy of 50% is adjusted to be 67%, representing the proportion of the maximum possible accuracy. In practice, the probability of intelligibility is a floating point probability in  $[0, 1]$ , which is typically compared to a threshold, the estimated probability of other words in the same phrase, or both, so the accuracy with which we can predict boolean intelligibility is only used as a benchmark by which we can measure the relative utility of different prediction methods.

## 7. Determining optimal feedback

We use the modeled probability of intelligibility of each word in a prompt word or phrase to help students improve their pronunciation by providing audiovisual feedback indicating which word(s) were pronounced the worst. How many words to indicate were not pronounced well after each utterance is an open question.

For words which are not considered sufficiently intelligible, we can use the DNN models to determine which identical numerical improvement to each phoneme's non-duration features improves the probability of word intelligibility the most. We can also see how increasing and decreasing each phoneme's duration improves the intelligibility of the word. Such adjustments to the features derived from automatic speech recognition may be more useful as products than sums to identify the specific phoneme(s) most in need of improvement in the less unintelligible word(s).



**Figure 3**

Determining feedback: Adjusting the feature scores for each phoneme changes the probability of intelligibility of the whole word. The adjustments which make the best changes signal which phoneme(s) need the most improvement.

## 8. Conclusion

Using PocketSphinx automatic speech recognition with improved phonetic accuracy features training DNN prediction models can help CAPT systems provide better intelligibility remediation. Researchers and commercial software publishers should try to understand the underlying reasons this technique is superior to the state of the art, and adopt it for improved CAPT outcomes.

## Acknowledgments

We thank 17zuoye.com of Beijing, China, Professor Seiichi Nakagawa, the Google Open Source Programs Office, and the Wikimedia Foundation for their kind financial support, suggestions, personnel resources, and educational infrastructure.

## References

- Australian Associated Press. 2017. Computer says no: Irish vet fails oral English test needed to stay in Australia. *The Guardian*. August 8.
- Chen, N.F. and Haizhou Li. 2016. Computer-assisted pronunciation training: From pronunciation scoring towards spoken language learning. In *Proceedings of the Asian-Pacific Signal and Information Processing Association (APSIPA) Annual Summit and Conference*, Jeju, South Korea. Review. [http://www.apsipa.org/proceedings\\_2016/HTML/paper2016/227.pdf](http://www.apsipa.org/proceedings_2016/HTML/paper2016/227.pdf).
- Dudy, S., S. Bedrick, M. Asgari, and A. Kain. 2017. Automatic analysis of pronunciations for children with speech sound disorders. Unpublished. <http://www.shirandudy.com/papers/automatic.pdf>.
- Ferrier, T. 2017. Australian ex-news reader with English degree fails robot's English test. *Sydney Morning Herald*. August 9.
- Huggins-Daines, D., M. Kumar, A. Chan, A.W. Black, M. Ravishankar, and A.I. Rudnicky. 2006. Pocketsphinx: A free, real-time continuous speech recognition system for hand-held devices. In *Proceedings of the IEEE International Conference on Acoustics, Speech and Signal Processing*, Toulouse, France. <https://www.cs.cmu.edu/~awb/papers/ICASSP2006/0100185.pdf>.
- Kibishi, H., K. Hirabayashi, and S. Nakagawa. 2015. A statistical method of evaluating the pronunciation proficiency/intelligibility of English presentations by Japanese speakers. *ReCALL*, 27(1):58–83. [http://www.slp.ics.tut.ac.jp/Material\\_for\\_Our\\_Studies/Papers/shiryou\\_last/e2014-Paper-01.pdf](http://www.slp.ics.tut.ac.jp/Material_for_Our_Studies/Papers/shiryou_last/e2014-Paper-01.pdf).
- Kibishi, H. and S. Nakagawa. 2011. New feature parameters for pronunciation evaluation in English presentations at international conferences. In *Proceedings of the Twelfth Annual Conference of the International Speech Communication Association (INTERSPEECH)*, pages 1149–52, Florence, Italy. [http://www.slp.ics.tut.ac.jp/shiryou/extra\\_2011/2011-IC-03.pdf](http://www.slp.ics.tut.ac.jp/shiryou/extra_2011/2011-IC-03.pdf).
- Llisterri, J. 2016. Computer-assisted pronunciation teaching bibliography. Archived at <http://j.mp/captbib>.
- Loukina, A., M. Lopez, K. Evanini, D. Suendermann-Oeft, A.V. Ivanov, and K. Zechner. 2015. Pronunciation accuracy and intelligibility of non-native speech. In *Proceedings of the Sixteenth INTERSPEECH*, Dresden, Germany. [http://www.aloukina.com/papers/Loukina\\_et\\_al.\\_intelligibility\\_Interspeech2015.pdf](http://www.aloukina.com/papers/Loukina_et_al._intelligibility_Interspeech2015.pdf).
- Qian, X., H. Meng, and F. Soong. 2015. A two-pass framework of mispronunciation detection and diagnosis for computer-aided pronunciation training. In *Proceedings of APSIPA Annual Summit and Conference*, Hong Kong, China. [http://www.apsipa.org/proceedings\\_2015/pdf/120.pdf](http://www.apsipa.org/proceedings_2015/pdf/120.pdf).
- Witt, S.M. and S.J. Young. 2000. Phone-level pronunciation scoring and assessment for interactive language learning. *Speech Communication*, 30(2):95–108. <http://svr-www.eng.cam.ac.uk/~sjy/papers/wiyo00.pdf>.