

# Pronunciation Assessment Proposal

**Jim Salsman**

**April 2019**

Item (1): Nakagawa et al's 2011 intelligibility vs. SRI's 1995 "Goodness of Pronunciation" assessments:

**82% agreement** with the accuracy of crowdworkers' transcriptions, up from 75% reported by the inventor

— [arxiv.org/abs/1709.01713](https://arxiv.org/abs/1709.01713)

## **State of the art:**

Educational Testing Service's SpeechRater 5.0

"System-Human agreement" is **58.4%** (Chen *et al.*, 2018.)

# Terms

**Speech recognition**

**Pronunciation assessment**

**Fluency assessment**

**(Usually requires accent and dialect adaptation)**

**Intelligibility prediction (Nakagawa et al., 2011)**

**Remediation**

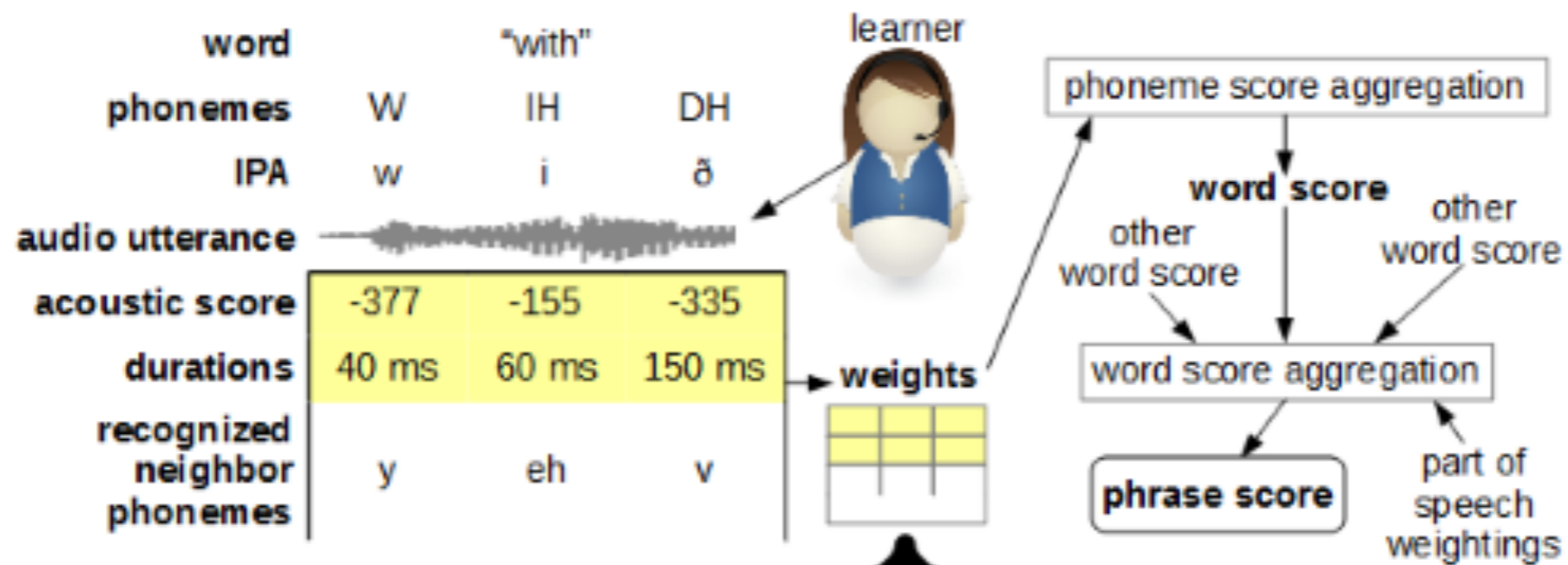
**Feedback -- natural speech or visual?**

**Interactions -- web app or download delay?**

**Sequencing -- good learner analytics?**

# Motivation

## Intelligibility assessment



Text-independent scoring

Learner recording exemplars

Native speaker exemplars

Authentic intelligibility remediation



no weighting of scores

Weight scores by many other students' attempts at pronunciation.

Weight scores by many known good pronunciation recordings.

Weight scores by their correlations with utterances yielding successful transcriptions.

**(A)**  
**worst**

**(B)**  
**poor**

**(C)**  
**better**

**(D)**  
**best**

# Why accent adaptation?



**14 accent loci in UK and Ireland,**

**3-4 in the USA and Canada,**

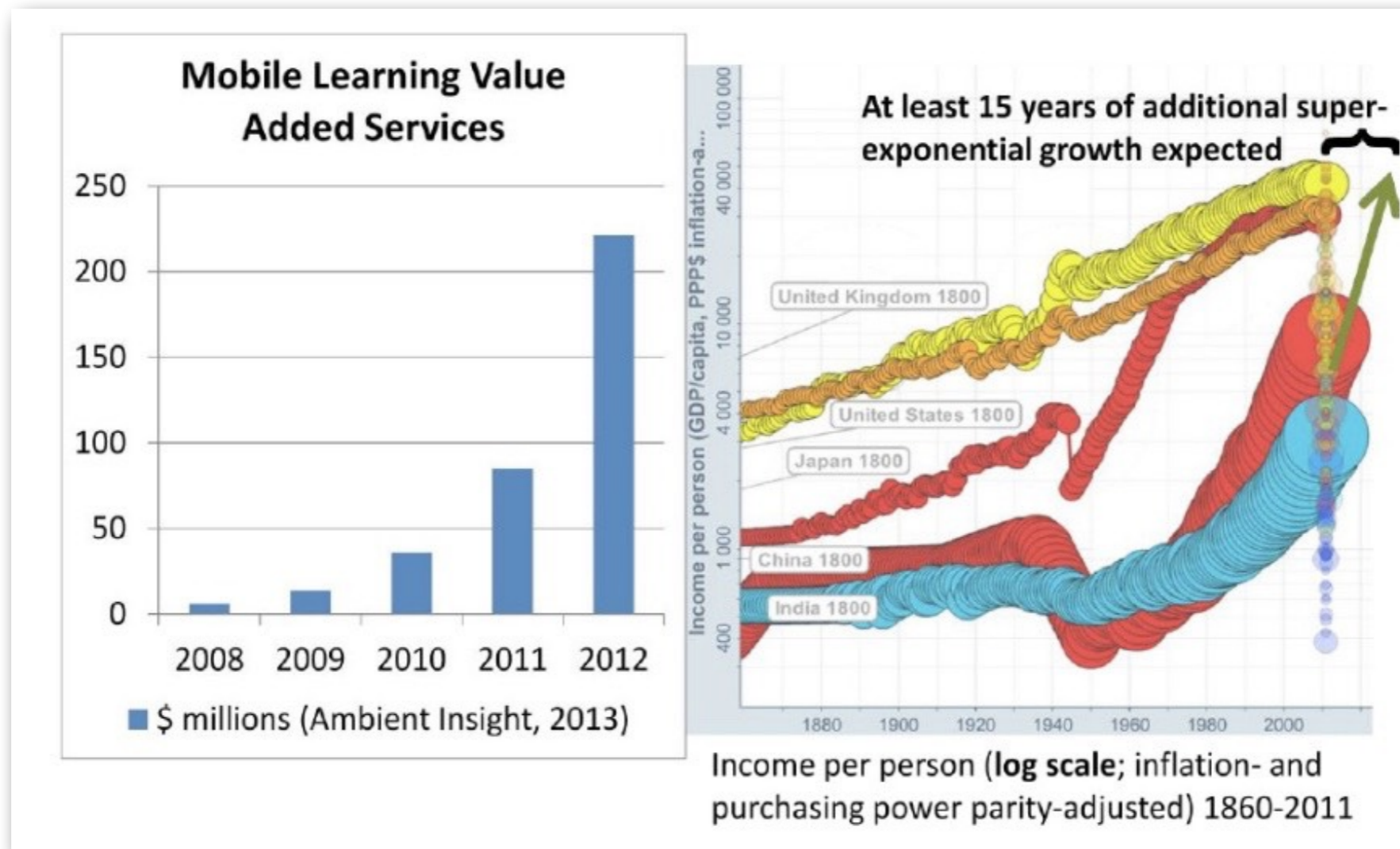
**Australia, South Africa, New Zealand, and**

**English as a Foreign Language everywhere**

**At least 150 feasibly discernible accents.**

**Fig. 1.** Accents of the British Isles.

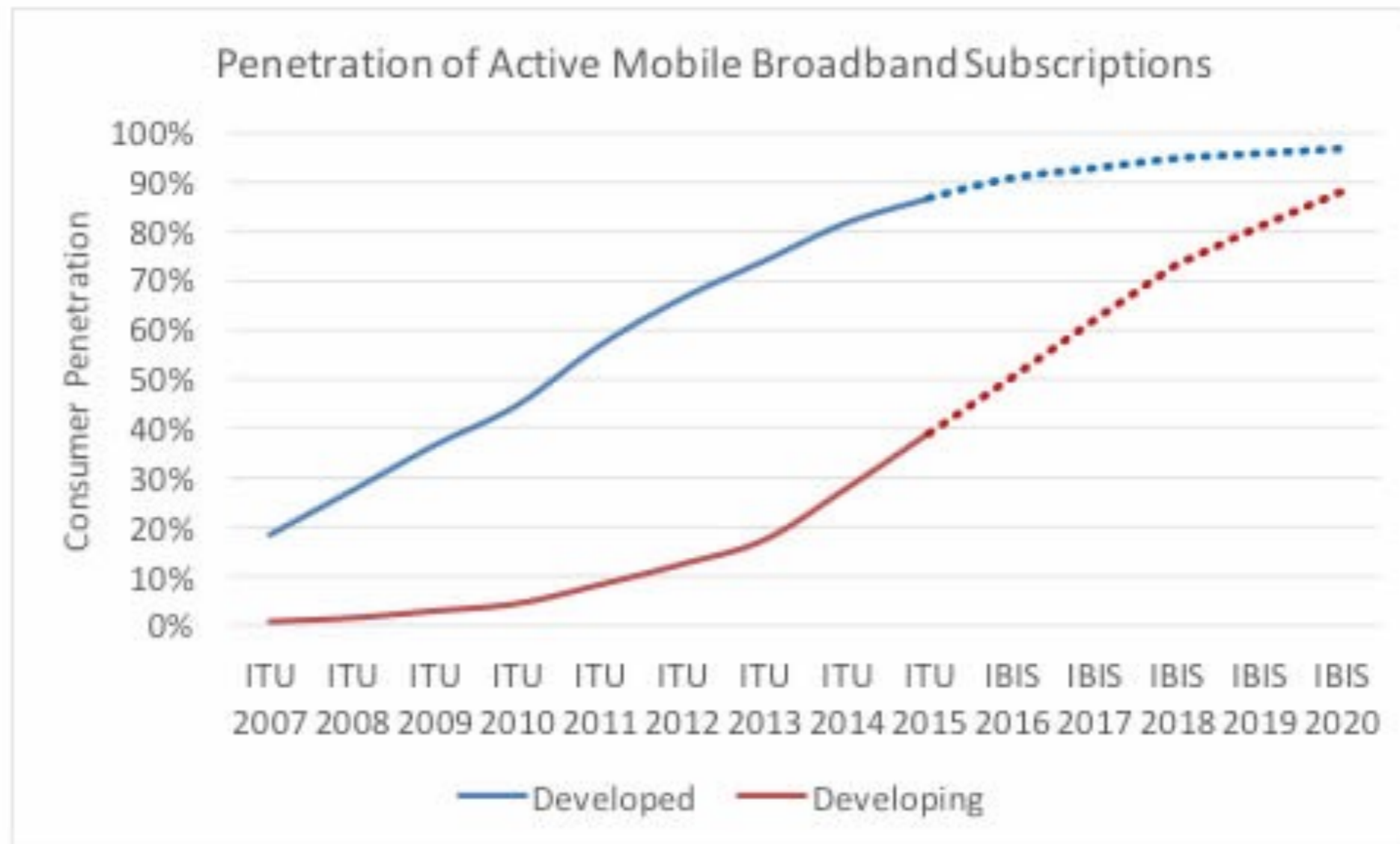
# Market size and opportunity



- Online language learning was a \$6 billion market in 2011. The global market for digital English language learning products reached \$1.8 billion in 2013, and exceeded \$3 billion in 2018.

# Market

## English language instruction



# Market size and opportunity



- The year-over-year growth rate was 20% in 2018. *WiseGuyReports* projects the global digital English language learning market will surpass \$18 billion by 2022.

# Goal

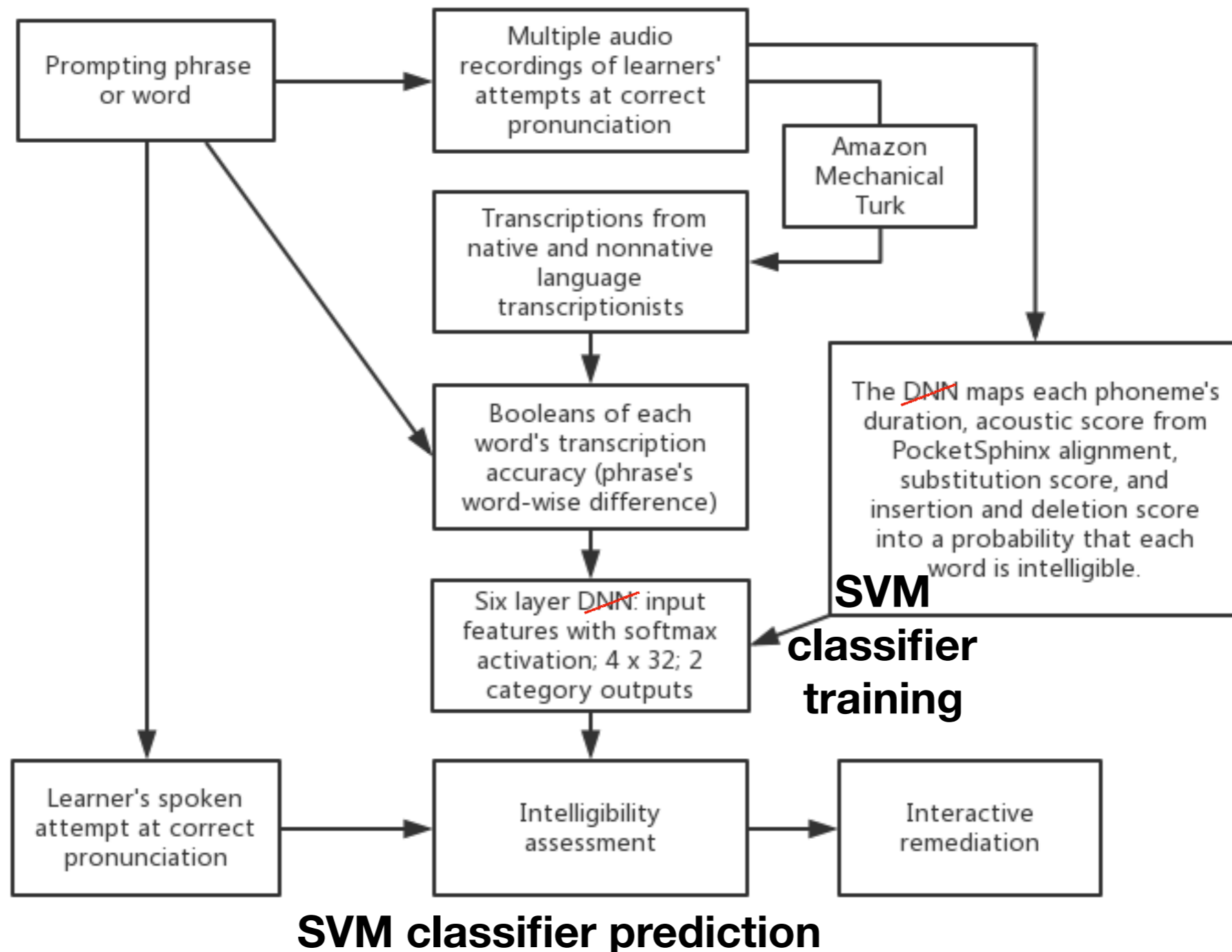
**The goal is to ask people to try to pronounce words, and some day phrases, in a way that speech recognition features predict will be correctly transcribed by those who hear the audio utterance. This technique can correctly adapt to spoken accents like vowel shifts, but not dialect.**

**Learners provide needed transcriptions of student (peer) speech.**



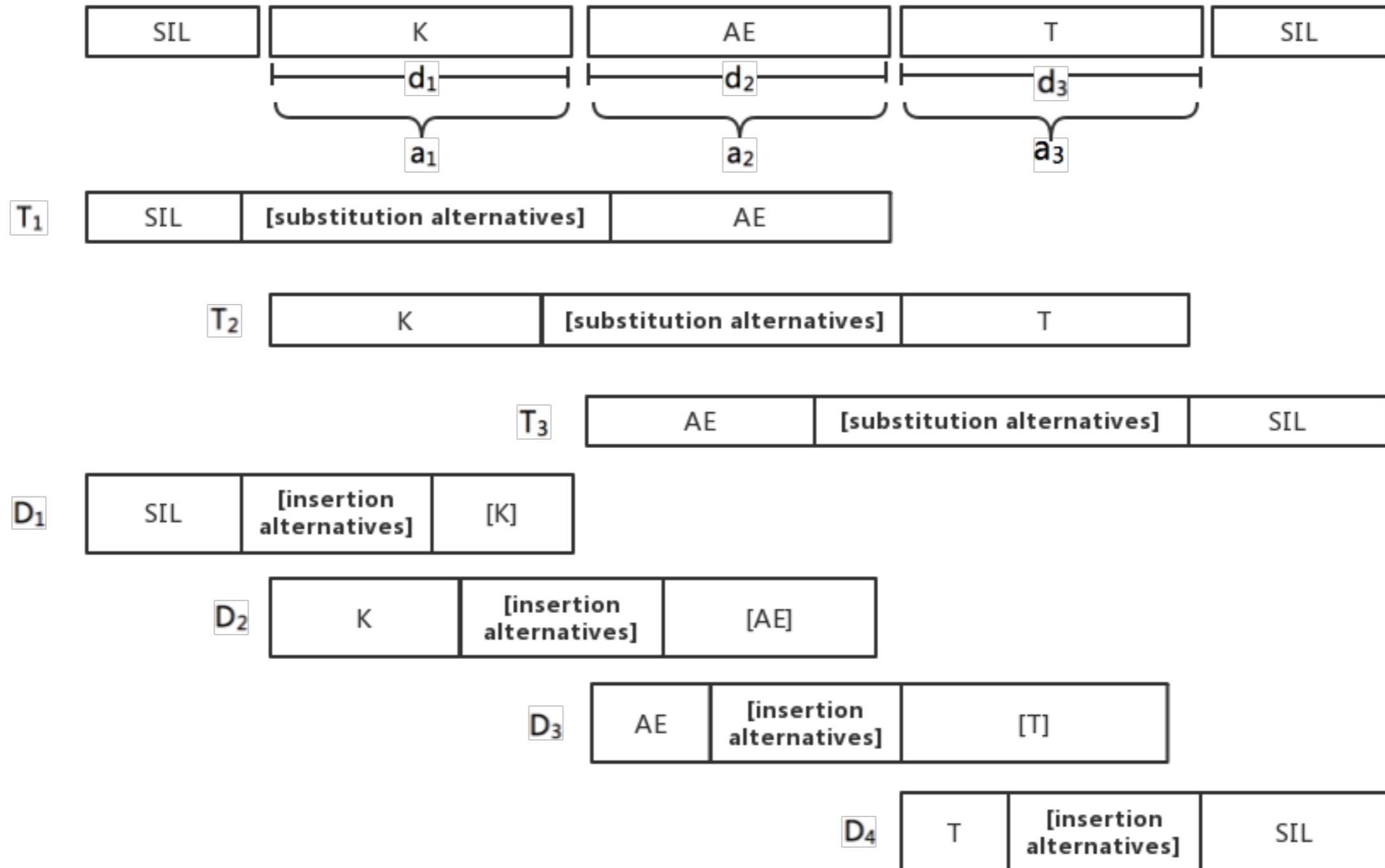
# Intelligibility remediation

## Data flow diagram



# Intelligibility remediation

## Scoring: 4 features/phoneme

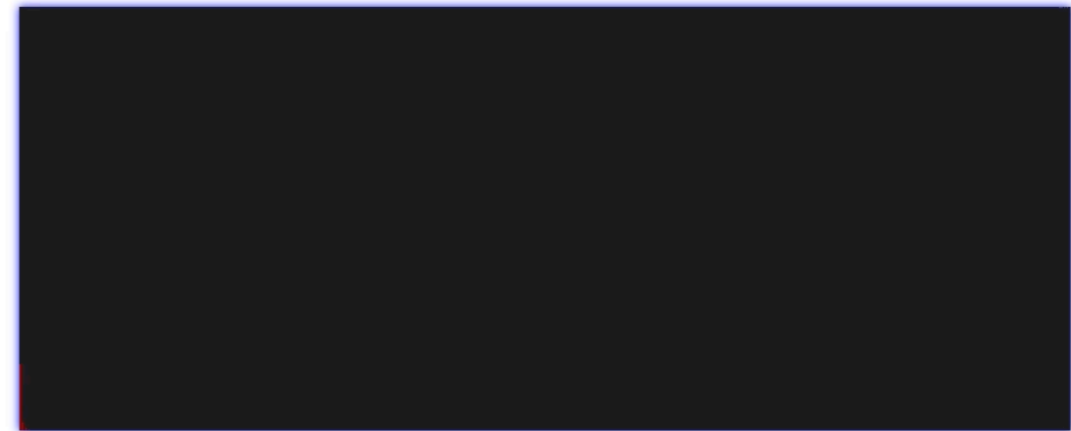




# Solution

- Natural spoken remediation feedback produces authentic skill improvements without the distraction of visual feedback.

Please pronounce "happy"



Record

Play

Evaluate

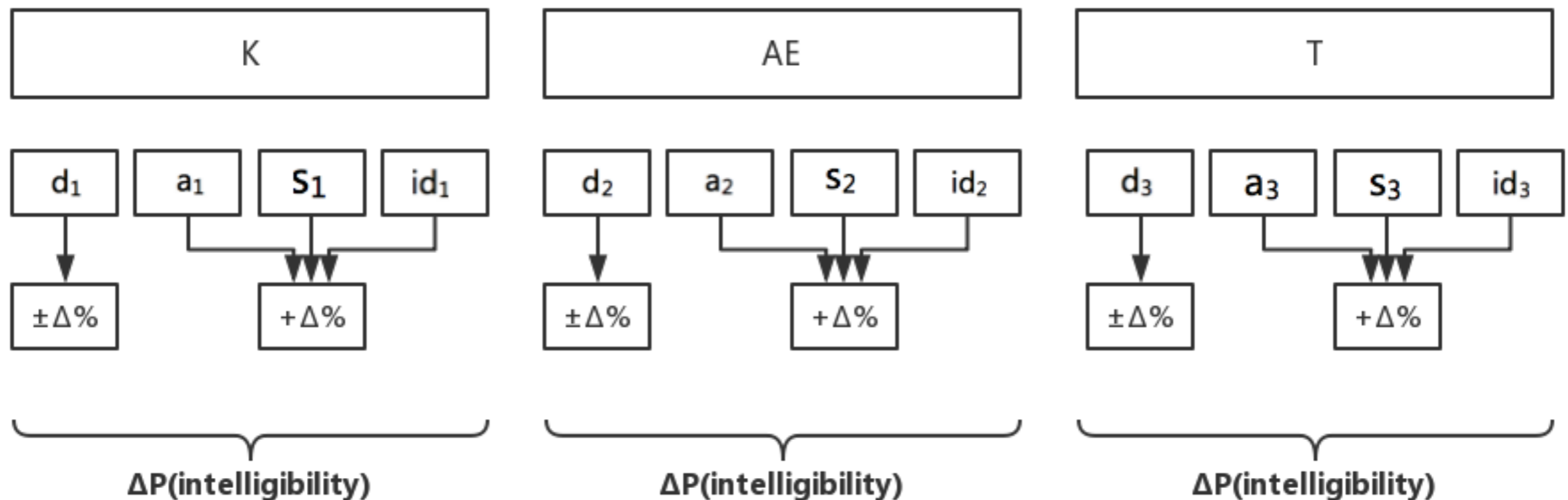
Email:

[jim@speakclearly.ic](mailto:jim@speakclearly.ic)

Exemplary

[Demo on YouTube](#)

## (2) Diphone with most room for improvement: "worst" phoneme(s)



The additional four vocal tract articulation features are all set to zero, and the ninth, proportion of neighbors less likely, is set to 1.0 for this step.

# Intelligibility remediation Manifest and plans

**Speech collection**

**Transcript collection**

**Transcript integration**

**Balancing**

**Sufficiency**

**Scoring**

**Feature extraction**

**SVM classifier**

**Phoneme with-most-room-for-improvement isolation**

**ID (email) -- adaptivity and payment processing integration: TBD**

**Exemplary flag for data collection: DONE**

**-- collecting transcripts from learners: IN PROGRESS**

**Multiple choice support tool done but not used yet**



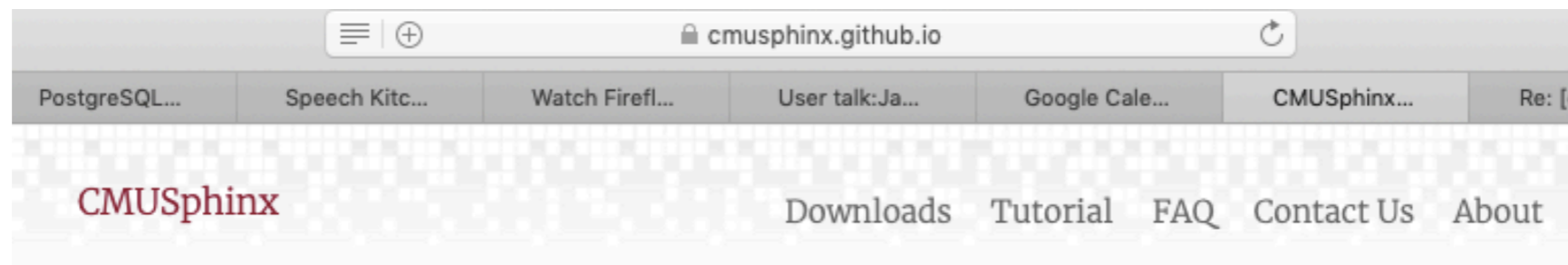
# Web pages

## API utilization in Javascript

```
194 mp3blob = new Blob(buffer, {type: 'audio/mp3'});
195
196 (window.XMLHttpRequest) ? req = new XMLHttpRequest() :
197   (window.ActiveXObject) ? req = new ActiveXObject("Microsoft.XMLHTTP") :
198   req = false; // cross platform for IE7 or something
199
200 req.open("post", "/");
201
202 formdata = new FormData();
203 formdata.append('url', '');
204 formdata.append('filetype', 'mp3');
205 formdata.append('word', document.getElementById("word").innerHTML);
206 formdata.append('email', document.getElementById('email').value);
207 formdata.append('exemp', document.getElementById('exemp').checked);
208 formdata.append('audio', mp3blob);
209 formdata.append('json', 'checked')
210
211 req.addEventListener('load', function(event) {
212   document.getElementById('upload').disabled = true;
213   var resp = JSON.parse(req.responseText);
214   var fUrl = '/rec/index.html?email='
215     + escape(document.getElementById('email').value);
216   if (resp.prob_good < 50 && resp.feedback != "") {
217     fUrl += '&word=' + escape(word) + '&feedback=' + escape(resp.feedback);
218   }
219   alert('Data sent, parsed response:\n' + req.responseText
220     + 'Visiting: ' + fUrl);
221   document.location = fUrl;
222 });
223 req.addEventListener('error', function(event) {
224   alert('Unable to upload.');
```

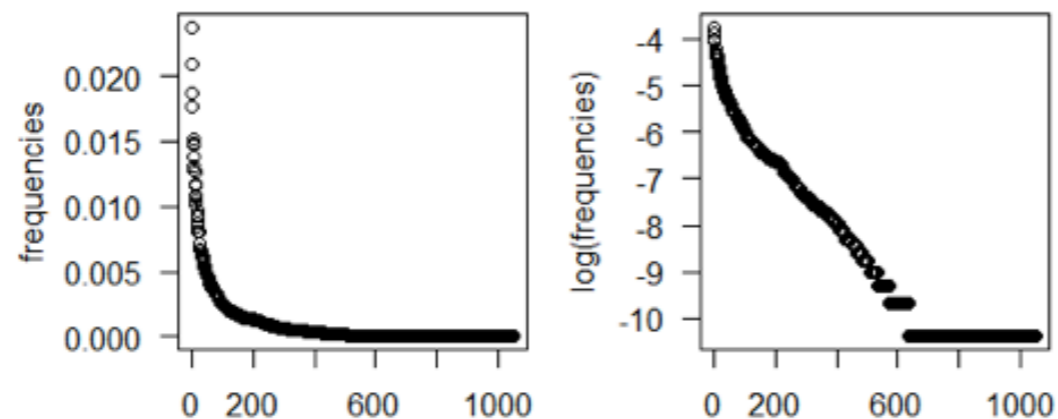


# (3) ~650 non-diphthong diphones for speech skill learner analytics, instead of phonemes or words



A *diphone* is the last part of one phoneme followed by the first part of another. Either phoneme could be silence, and they can be the same phoneme. Diphthongs include diphones in them.

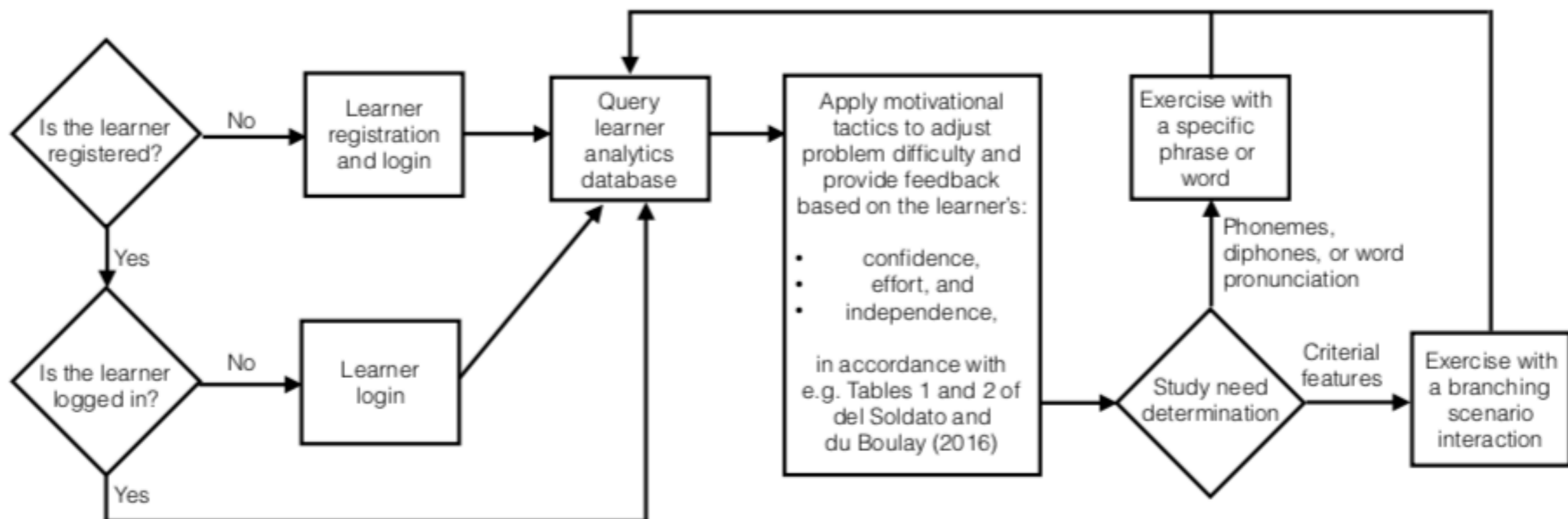
This list of the top 4,800 words by frequency in English speech was used with CMUDICT to create the following list of 1,052 diphones without diphthongs by approximate prevalence.



UH\_R 2.376%, AH\_N 2.083%, T\_SIL 1.863%, Z\_SIL 1.758%, SIL\_S 1.514%, IY\_SIL 1.486%,  
EH\_IY 1.465%, D\_SIL 1.387%, N\_SIL 1.313%, S\_SIL 1.270%, SIL\_K 1.264%, R\_SIL 1.168%,  
AA\_IY 1.156%, SIL\_P 1.078%, IH\_NG 1.063%, S\_T 1.048%, AH\_L 1.017%, NG\_SIL 0.973%,

# Adaptivity

Learner analytics sequencing and branching scenario transitions



Reference: du Boulay, B. and del Soldato, T. (2016) "Implementation of Motivational Tactics in Tutoring Systems: 20 years on," *International Journal of Artificial Intelligence in Education*, 26(1):170-182, <http://users.sussex.ac.uk/~bend/papers/motivation-revised2.pdf>

# (4) JavaScript PocketSphinx.js recognition on the client web browser

```
load_prompts_utt.py  pocketsphinxjs-config.h.in  psRecognizer.cpp  load_words.py  combo-new.dict

1  #include "psRecognizer.h"
2  #include "pocketsphinxjs-config.h"
3
4
5  namespace pocketsphinxjs {
6      typedef std::map<std::string, std::string> StringsMapType;
7      typedef std::map<std::string, std::string>::iterator StringsMapIterator;
8
9      // Implemented later in this file
10     ReturnType parseStringList(const std::string &, StringsSetType*, std::string*);
11
12     Recognizer::Recognizer(): is_fsg(true), is_recording(false), current_hyp(""),
13     • grammar_index(0) {
14         Config c;
15         if (init(c) != SUCCESS) cleanup();
16     }
17
18     Recognizer::Recognizer(const Config& config) : is_fsg(true), is_recording(false),
19     • current_hyp(""), grammar_index(0) {
20         if (init(config) != SUCCESS) cleanup();
21     }
22
23     ReturnType Recognizer::reInit(const Config& config) {
24         ReturnType r = init(config);
25         if (r != SUCCESS) cleanup();
26         return r;
27     }
28
29     ReturnType Recognizer::addWords(const std::vector<Word>& words) {
30         if (decoder == NULL) return BAD STATE;
```

# **(5) Data collection: Words, speech and transcripts**

**700 words (for comparison the Cambridge/EC English Profile has 6,500 words in levels A1-C2) and phrases;**

**30-60 recordings per word;**

**4-12 transcripts per recording; and**

**4 numeric features per phoneme, upgraded to 10.**

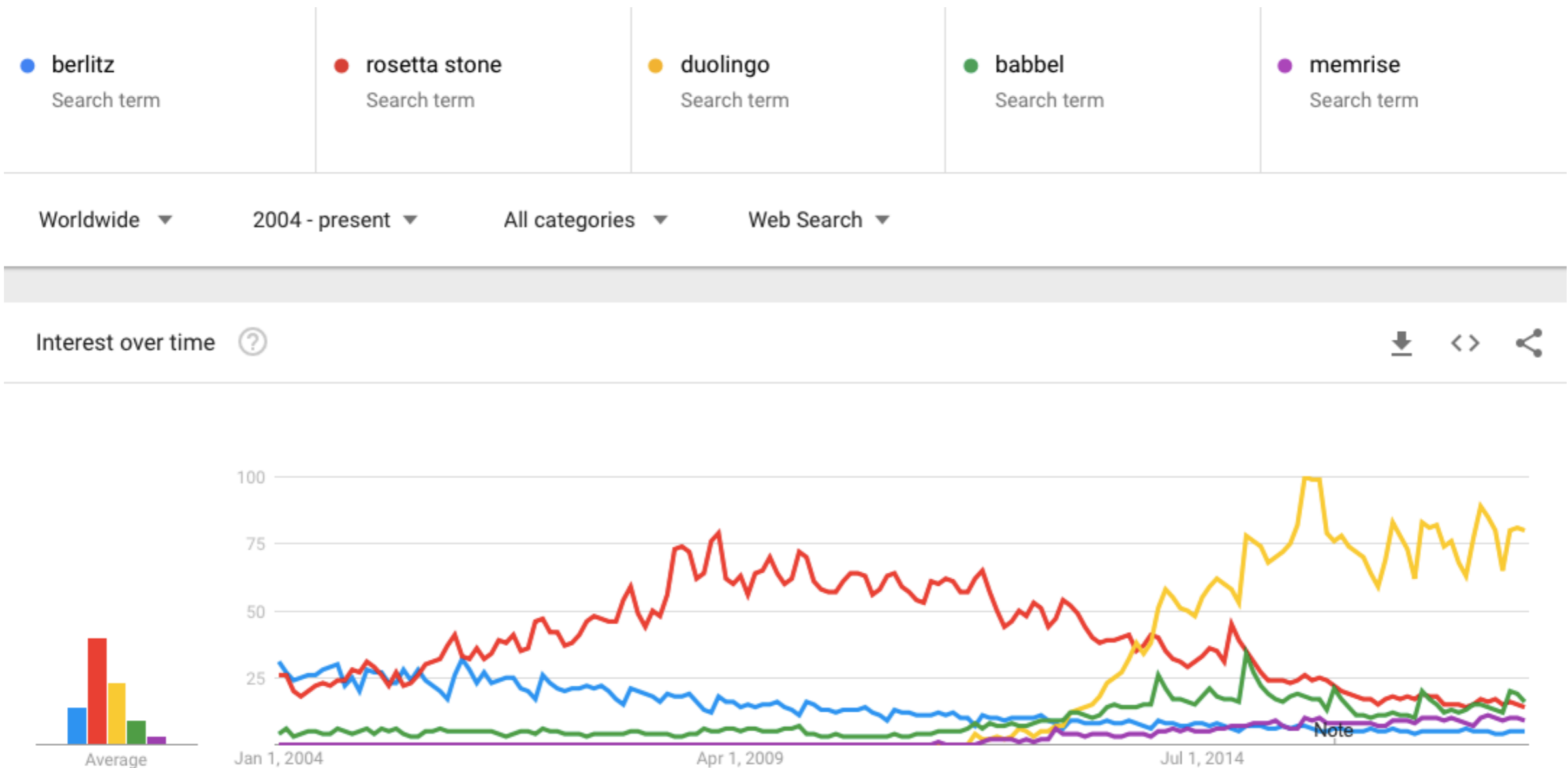
# Data collection Balancing

	Current	Goal	Balancing
<b>Prompts (word or phrase)</b>	700	7000	Vocabulary grade level (eg. A1, A2, B...)
<b>Recordings</b>	30 per prompt	60 per prompt	Requires both good, completely wrong, and marginal
<b>Transcripts</b>	4 per recording	8 per recording	Beware of corruption from lazy and other defectors
<b>Exemplary recordings</b>	15 per prompt (40 words)	s*4 per prompt? (2 gender x 2 age)	s needs to be large enough for balancing recordings

**Schema:**

**Users,  
Authenticators,  
Words,  
Utterances,  
Prompts,  
Topics,  
Choices,  
Lessons,  
Schools.**

# (6) Proposal TBD, e.g. third party restricted market sale free from ongoing cost center resource drains.



**Questions? Email:  
jim@talknicer.com**