

AI-assisted analysis of phonological variation in English

Special session on Deep Phonology, AMP 2025

UC Berkeley September 27, 2025

Virginia Partridge, Joe Pater, Parth Bhangla,

Ali Nirheche and Brandon Prickett

{vcpartridge, pater}@umass.edu

UMassAmherst



These slides, and references, are available at <https://websites.umass.edu/pater/handoutsslides/>

Background

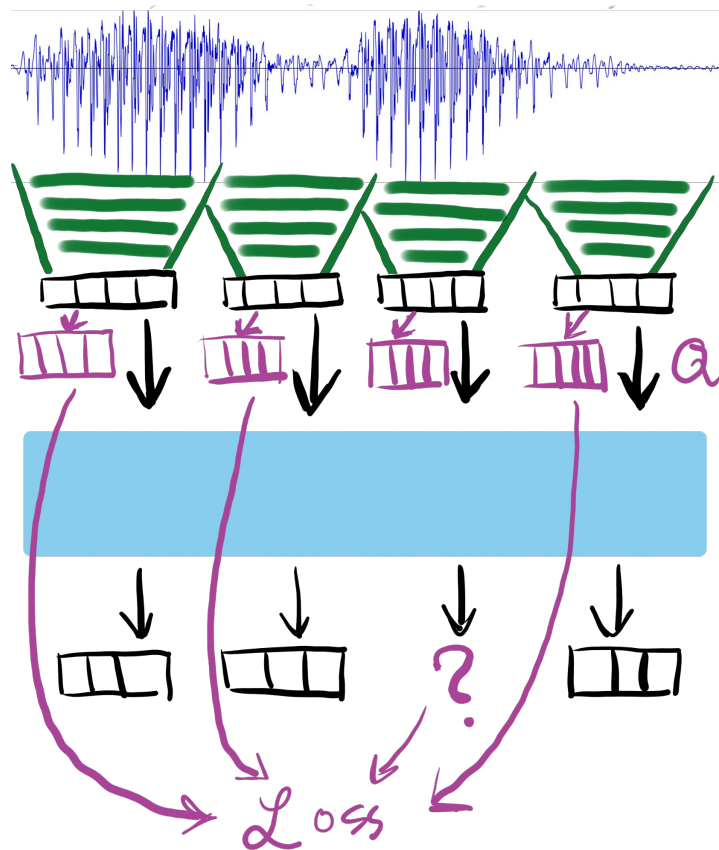
Despite the dramatic recent advances in speech recognition technology, automated phonetic transcription is not widely used in linguistics (as far as we know). This is especially surprising since the paper introducing Wav2Vec 2.0 ([Baevski et al. 2020](#)) presents SOTA results on a TIMIT phone recognition benchmark.

Our goal is to increase the usefulness of automated phonetic transcription for the study of phonological variation in English (variation within and across speakers and varieties)

In this paper, we:

- Situate our work within current research on automated transcription
- Present a Wav2Vec 2.0 model fine-tuned on Buckeye ([Pitt et al. 2005](#)), and compare it to other models, using TIMIT ([Garofolo et al. 1993](#)) as a test set
- Provide a web-based implementation, with Praat textgrid input and output
- Discuss next steps (and also issues with the standard TIMIT benchmark)

Wav2Vec 2.0 Pre-training: Sounds to vectors



Convolutional Layers: Snapshots of what's happening in audio at a particular point in time

Transformer Layer: Take into account relative positions of sounds in context

Masks + Quantization + Loss Function: Perform *self-supervised learning* by hiding sound representations that must be predicted using a probabilistic function

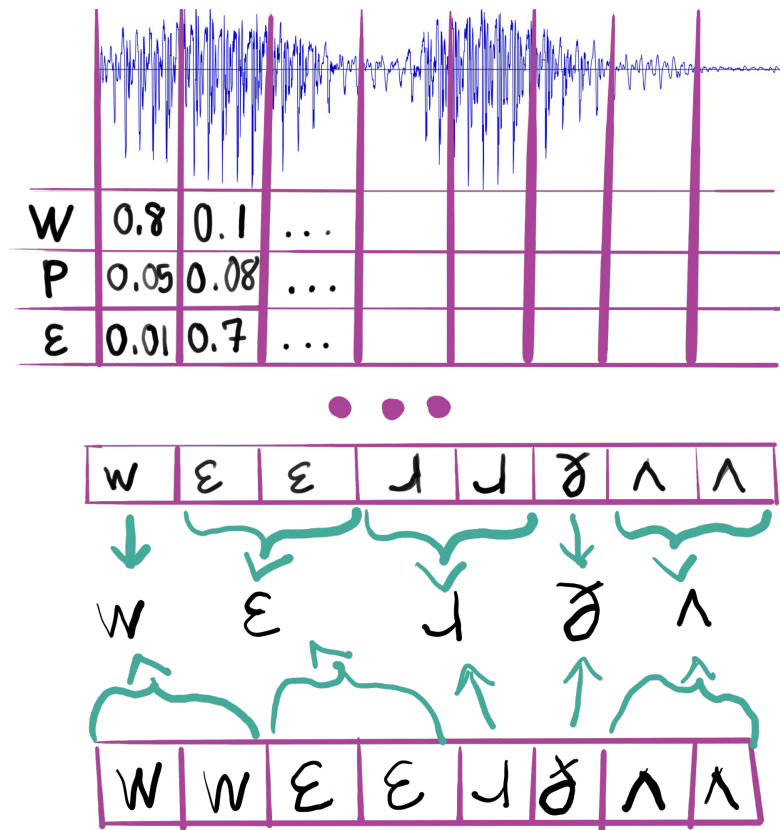
Facebook Research has open sourced many [many pre-trained models](#) for others to use.

Fine-tune: Add phone prediction

Bring your own audio with transcriptions.

Using Connectionist Temporal Classification loss, **maximize the probability of the correct output** by:

- 1) Predicting each symbol's **probability at each time segment**
- 2) **Merging** repeated symbols across neighboring time segments
- 3) Maximize the marginal probability of **alignments** that lead to the correct output



Multilingual automated transcription

The primary focus of contemporary research on automated phonetic transcription is the development of universal models, trained on data from multiple languages, and applicable in principle to any language.

This research includes fine-tuned Wav2Vec 2.0 models ([Xu et al. 2021](#), [Taguchi et al. 2023](#)) as well as the quite different Allosaurus ([Li et al. 2020](#)), which emphasizes the distinction between phoneme and phone recognition.

The latest contribution to this line of research is [Zhu et al. \(2025\)](#). It advances the SOTA on benchmarks, but also recognizes that “error analysis reveals persistent limitations in modeling socio-phonetic diversity, underscoring challenges for future research”

Zhu et al. illustrate these limitations by testing the model on the Buckeye English corpus. We'll be doing this for multilingual models as well, but not with theirs, since we were unfortunately unable to get it to run.

Grapheme-to-phoneme vs. actual transcriptions

Multilingual transcription leverages the availability of large amounts of orthographically transcribed speech by applying grapheme-to-phoneme conversion (e.g. Epitran; [Mortensen et al. 2018](#)) to obtain phonetic transcriptions.

These are similar to dictionary transcriptions, and are of course not equivalent to transcriptions of the individual utterances:

- The pronunciation of individual words can vary across speakers, across utterances by an individual speaker, and across phonetic and phonological contexts
- The pronunciation of phonemes can differ within words. Allosaurus also uses phoneme-to-allophone conversion, but the transcriptions are still relatively abstract

AutoIPA

We call our model “AutoIPA” (AI-assisted in our title is thus ambiguous)

It is a Wav2Vec 2.0 pre-trained model (facebook/wav2vec2-large-xlsr-53) fine-tuned on the Buckeye corpus:

- 40 speakers, balanced for gender and age (over 40 vs. under 40)
- White residents of the greater Columbus area (Northern Midland)
- 30 to 60 minutes of conversational speech for each speaker - about 20 hours total
- Phonemic transcription + vowel nasalization, flap and glottal/glottalized stop (we adopt the [Seyfarth and Garellek 2020](#) revisions), includes syllabic sonorants
- No distinction between stressed/unstressed $\Lambda/\text{ə}$ or $\text{ɜ}/\text{ə}$ (“perceived quality alone” [Kiesling et al. 2006: 18](#); see relatedly [Lindsay 2022](#)).

As far as we know, this is the first time Buckeye has been used in the training of an automated transcriber, though Buckeye and TIMIT have been used in training phone alignment models (e.g. [Kreuk et al. 2020](#)), and Zhu et al. use it as a test set.

Processing the Buckeye data

Data from 24 speakers used for training, 8 for development, and 8 for testing

- Each demographic is represented in equal proportion across these data sets
- E.g., 6 younger women, 6 older women, 6 younger men, and 6 older men in training
- Success on test data requires generalization to new speakers

Training data divided into samples 8-12 seconds in length for Wav2Vec 2.0 fine-tuning

ARPABET transcriptions converted to IPA; the symbol set with our ARPABET to IPA mappings can be seen [here](#). We used Phonecodes ([Hasegawa-Johnson 2019](#)): our [GitHub fork](#), our [Python package](#)

The task is phone (rather than feature) prediction, affricates and diphthongs treated as single phones, and nasalization and syllabicity are properties of individual phones

Evaluation Metrics: Phone Error Rate (PER)

What's the **ratio of phones** we get wrong in each test sample?

- 1) Compute the edit distance between your prediction p and your reference r to find the total number of phone insertions, substitutions or deletions between them.
- 2) Normalize by the number of phones in the reference r .

$$per(p, r) = \frac{edit_distance(p, r)}{length(r)}$$

$per([bop], [pop]) = \frac{1}{3}$

$per([bop], [po]) = 1$

$per([bip], [po]) = 1.5$

LOWER IS BETTER!

Evaluation Metrics: Phone Feature Error Rate (PFER)*

How many **articulatory features** are we getting wrong?

- 1) Compute edit distance to find phone substitutions, deletions and insertions between your prediction p and your reference r .
- 2) Sum up the total cost of edit distance errors between aligned phones:
 - Phone deletions and insertions each cost 1
 - For substitutions, each feature mismatch costs $1/24$ (there are 24 features in the table). This is the normalized Hamming distance between the articulatory features of the phones.

$$\text{pfer}([\text{b}\text{o}], [\text{p}\text{o}]) = 0.04166$$

$$\text{pfer}([\text{b}\text{o}\text{p}], [\text{p}\text{o}]) = 1.04166$$

$$\text{pfer}([\text{p}\text{ɛ}], [\text{p}\text{o}]) = 0.125$$

LOWER IS BETTER!

* Misleadingly **not really a rate**, but does obey the triangle inequality

Evaluation Metrics: Implementation

We use the [PanPhon Python library](#) ([Mortensen et al. 2016](#)) to calculate phone and feature edit distances.

For ease of use and reproducibility, we've published [a wrapper around PanPhon's distance metrics](#) that is compatible with the [HuggingFace evaluate package](#). This makes them easy to compute when working with standard Python libraries for transformer-based speech recognition models.

The results we report for models here are average PER and PFER across samples in the test corpus.

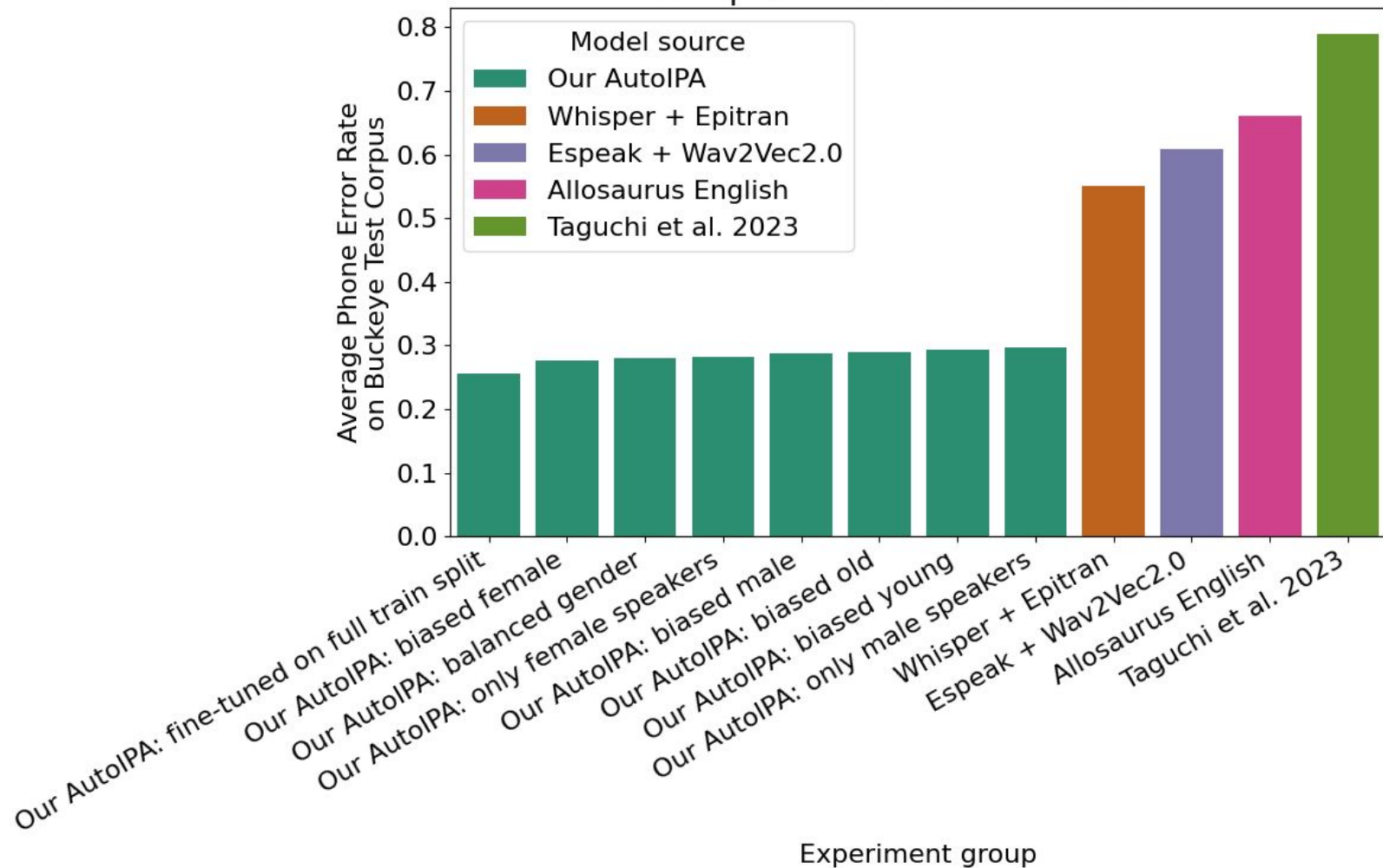
A caution on model comparison

It is completely expected that AutoIPA will outperform the multilingual models!

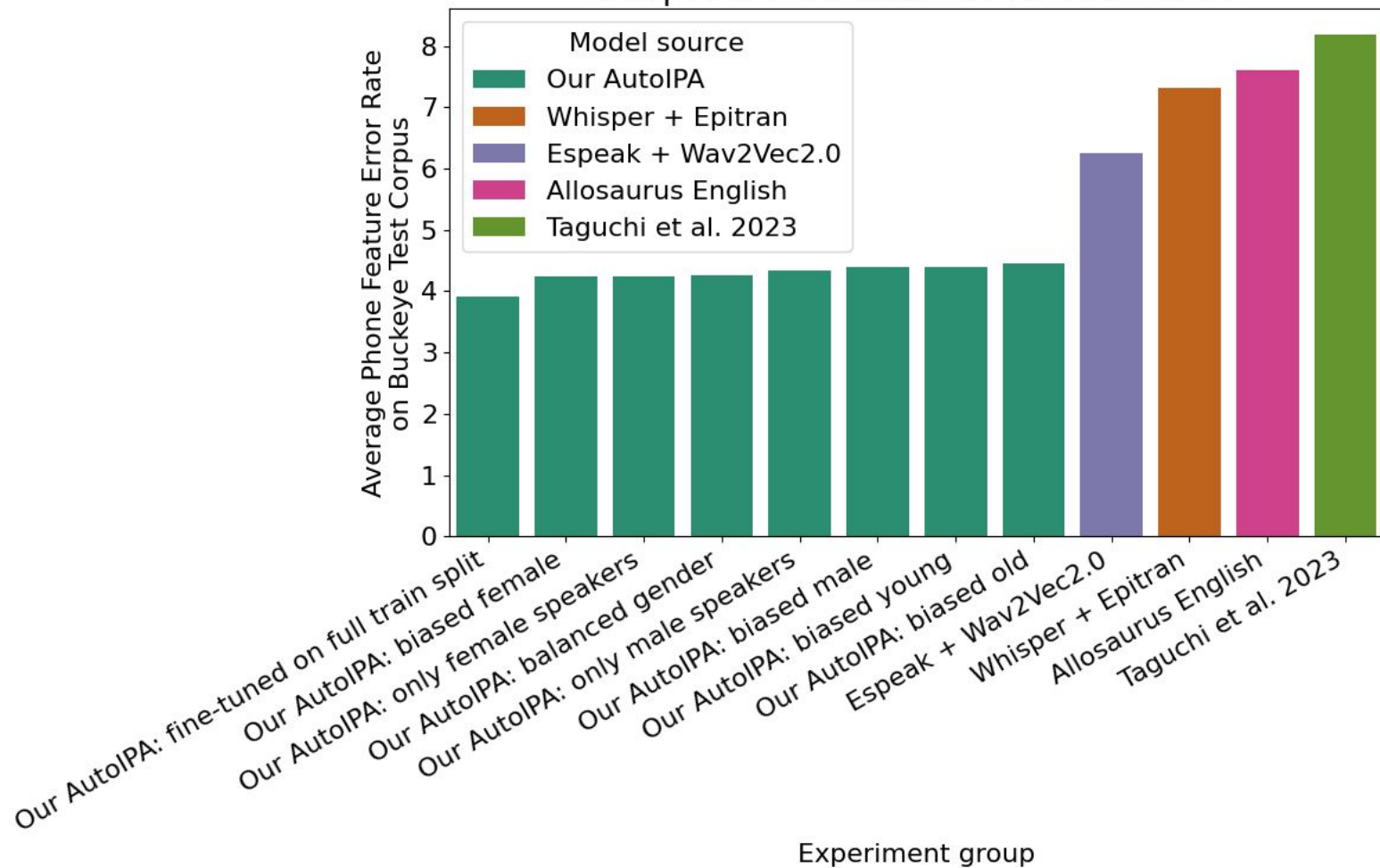
But we need *some* way of relativizing the performance of AutoIPA

Future work might take our results as a benchmark to improve on

Comparison of Phone Error Rates



Comparison of Phone Feature Error Rates



Some example transcriptions

The following are examples of transcriptions of utterances from the Buckeye test set by the models discussed here. The orthographic and phonetic transcriptions at the top are from Buckeye.

We have inserted spaces between words to increase readability.

s3802a_Utt33 (Older male)



where the trees were mature and they didn't really ruin our development and cut down
a lot of the large trees

wɛɪ ðɪ tʃɪz wɪ mætʃuɪ ɪn neɪ dɪŋ ɹɛli juːɹ̩ ɔɪ dɪvɛləpmɛnt ɪ kʌt daʊ̃r ɔ lɑːrɫv ðʌ lɑːdʒ tʃɪz

wɛɪ ðɛ tʃɪz wɪ mætʃuɪ ɛn neɪ dɪŋ ɹɛli juːɫn ɔɪ dɪvɛlpmɛ? ɪ kʌt daʊ̃r ʌ lɑːɫ ðʌ lɑːdʒ tʃɪz

AutoIPA

ðə tʃɪz wɪ mætʃuɪ ænd ðej dɪdntɪ ɹɪli juːən ɔwɪ dɪvɛləpmɛntæ kʌt daʊn ə lɑt ʌv ðə lɑːdʒ
tʃɪz

Allosaurus

wɛɪ ðə tʃɪz wɪ mætʃuɪ ænd ðej dɪdnt ɹɪl juːən ɔwɪ dɪvɛləpmɛnt ænd kʌt daʊn ə lɑt ʌv ðə
lɑːdʒ tʃɪz

Whisper + Epitran

or ðɛ trɪz wu mɪtʃɔj ɒn eɪ dɪn wɛɹvɪ wɛn ɔɪ dɪvɛwpmɛn kɛ daɪn ɔlɑːrɑ dɔ lɑdʒ trɪs

Taguchi

wɜː ðə tʃɪz wɜː mætʃuːɹ æn ðeɪ dɪdən ɹɪli juːɪn ɔɪ dɪvɛlpmɛntən kʌt daʊn ɐ lɑːrɫv ðə
lɑːdʒ tʃɪz

ESpeak+Wav2Vec 2.0

s3801b_Utt80 (Older male)



yknow they benefited a lot from it now

ji nou ðei benɛfɪtɪd ʌ lɔʔ fɪʌm ʌt naʊ

ji ŋou ðei beŋʌfɪtɪɾ ʌ lɔʔ fɪʌm ʌʔ naʊ AutoIPA

ðej benəfɪtɪd ə lɔt fɪʌm naʊ

Allosaurus

ðej benəfɪtɪd ə lɔt fɪʌm ɪt

Whisper + Epitran

i nɛv ðei benɛfedɪd ɛ aː fɛm nɛw

Taguchi

iː nə ðei benɪfɪtɪd ɐ lɔːt fɪʌm naʊ

ESpeak+Wav2Vec 2.0

s3902a_Utt109 (Younger female)



still just sitting down and my computer sits next to it VOCNOISE um

stɪl dʒʌs sɪtɪŋ daʊn ɪ maɪ kəmˈpjʊdə sɪts nekst tu ɪt ʌm

stɪl dʒʊs sɪtɪŋ daʊ̃ ʌm maɪ kəmˈpjɪrɪ sɪts nekst tu ʌt ʌm

AutoIPA

stɪl dʒʌst sɪtɪŋ daʊn ʌn maɪ kəmˈpjʊtɪ sɪts nekst tə ɪt

Allosaurus

stɪl dʒʌst sɪtɪŋ daʊn ʌn maɪ kəmˈpjʊtɪ sɪts nekst tə ɪt

Whisper + Epitran

stəʊ: dwɪ sɪˈerɪn daɪn ʌ maɪ kɪmpˈtʃɪrɪ sɪts nekst tə wɪ ʌm

Taguchi

stɪl dʒʌs sɪtɪŋ daʊn ʌ maɪ kəmˈpiːrə sɪtss nekstuːɪt ʌm

ESpeak+Wav2Vec 2.0

Top edit distance errors on the Buckeye test split

These are the **most common mistakes** we saw from a model fine-tuned on all data in the Buckeye train split (ginic/full_dataset_train_3_wav2vec2-large-xlsr-53-buckeye-ipa).

Substitution	Percentage of all substitutions
$\text{I} \rightarrow \Lambda$	5.27 %
$\text{I} \rightarrow \varepsilon$	4.70
$\Lambda \rightarrow \text{I}$	4.02
$\text{r} \sim \rightarrow \text{n}$	2.57
$\varepsilon \rightarrow \Lambda$	2.32
$\text{o} \rightarrow \Lambda$	2.29
$\text{i} \rightarrow \text{I}$	2.28

Deletion	Percentage of all deletions
I	14.74 %
v	9.34
Λ	9.09
I	5.88
n	5.70
t	4.95
l	3.64

Insertion	Percentage of all insertions
I	9.76 %
t	8.57
Λ	7.68
n	7.57
d	4.37
ε	4.13
v	3.95

Vowel Error Rates

For vowels v in the vocabulary V , we want to see **which are most challenging** across the entire corpus.

Which vowels is the model getting **wrong most frequently**?

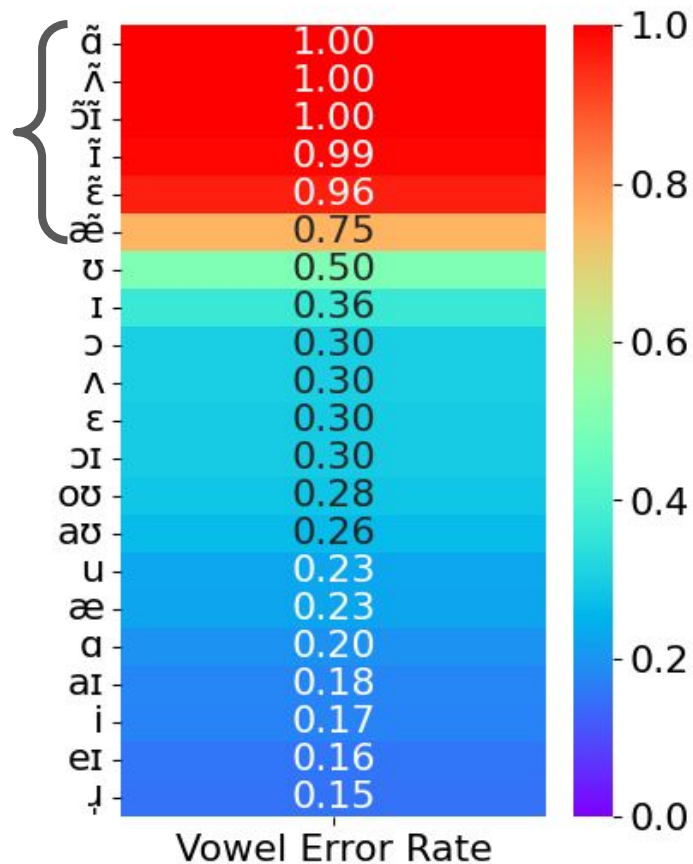
$$\textit{count_errors}(v) = \sum_{w \in V} \textit{count_substitutions}(v, w) + \textit{count_deletion}(v)$$

$$\textit{error_rate}(v) = \frac{\textit{count_errors}(v)}{\textit{total_count}(v)}$$

AutoIPA Buckeye Test Set Vowel Error Rates (Descending worst to best)

We don't do well on nasalized vowels.

They are relatively rare overall, which makes them more challenging for the model to recognize. They are also apparently not well transcribed by Buckeye's humans.



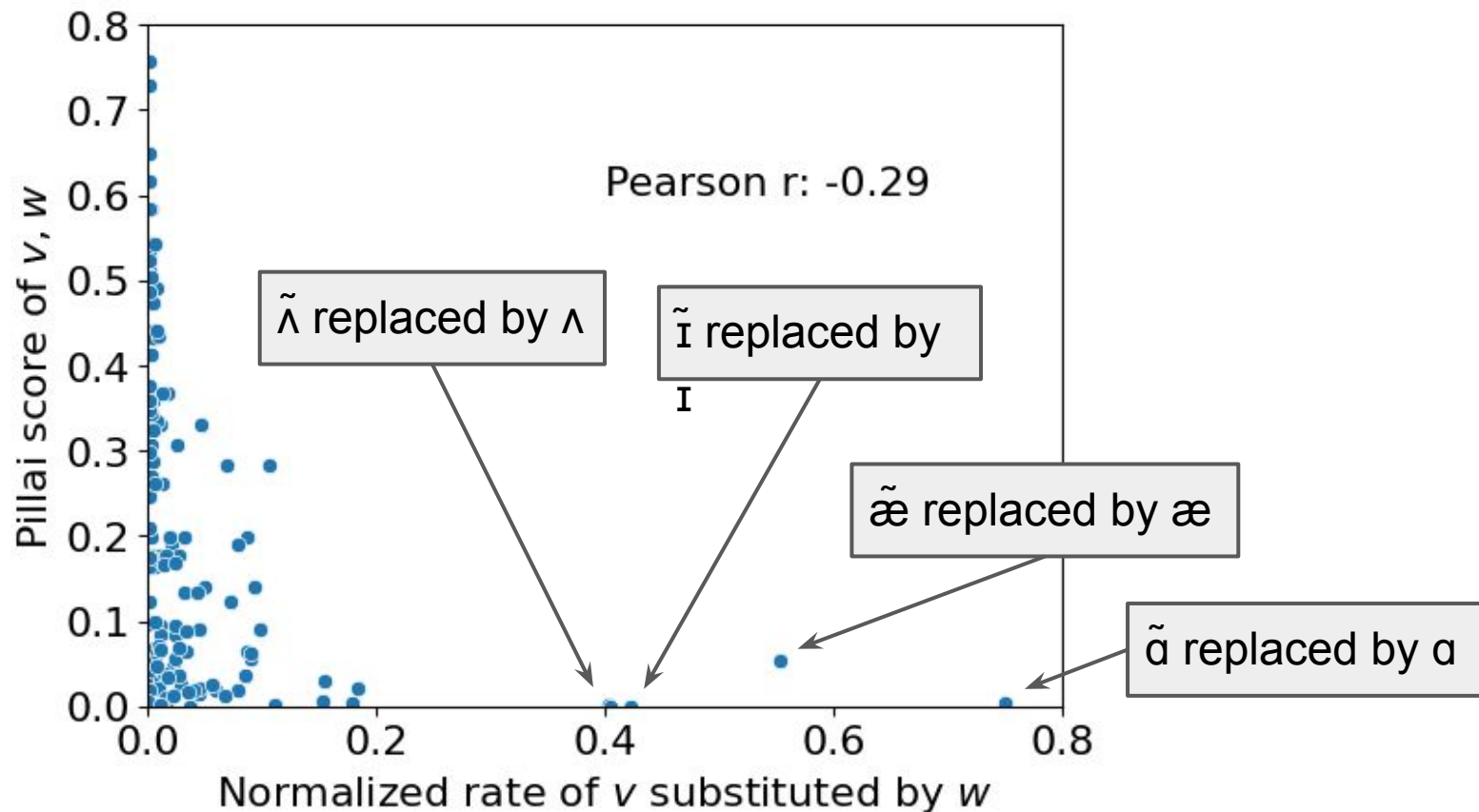
Pillai scores

Pillai scores are a measure of distributional overlap used in sociolinguistics to quantify the degree of vowel merger ([Hay et al. 2006](#), [Stanley and Sneller 2023](#)).

We calculated Pillai scores based on F1 and F2 measurements of the Buckeye (and TIMIT) vowels to see if they correlated with error rates

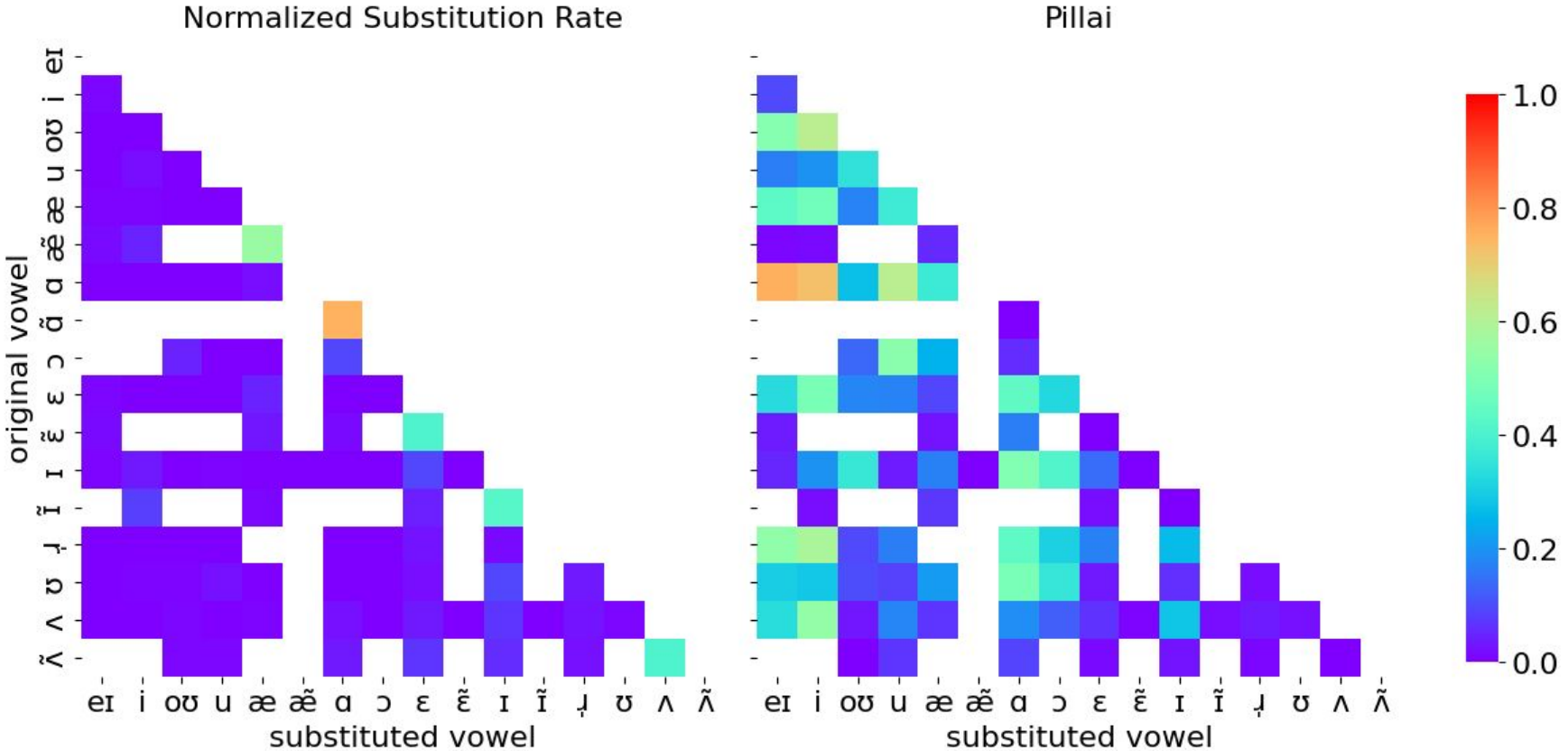
A high Pillai score indicates a low degree of overlap, or a high degree of contrast

Relationship between Pillai score within Buckeye and substitution rate in the Buckeye test split



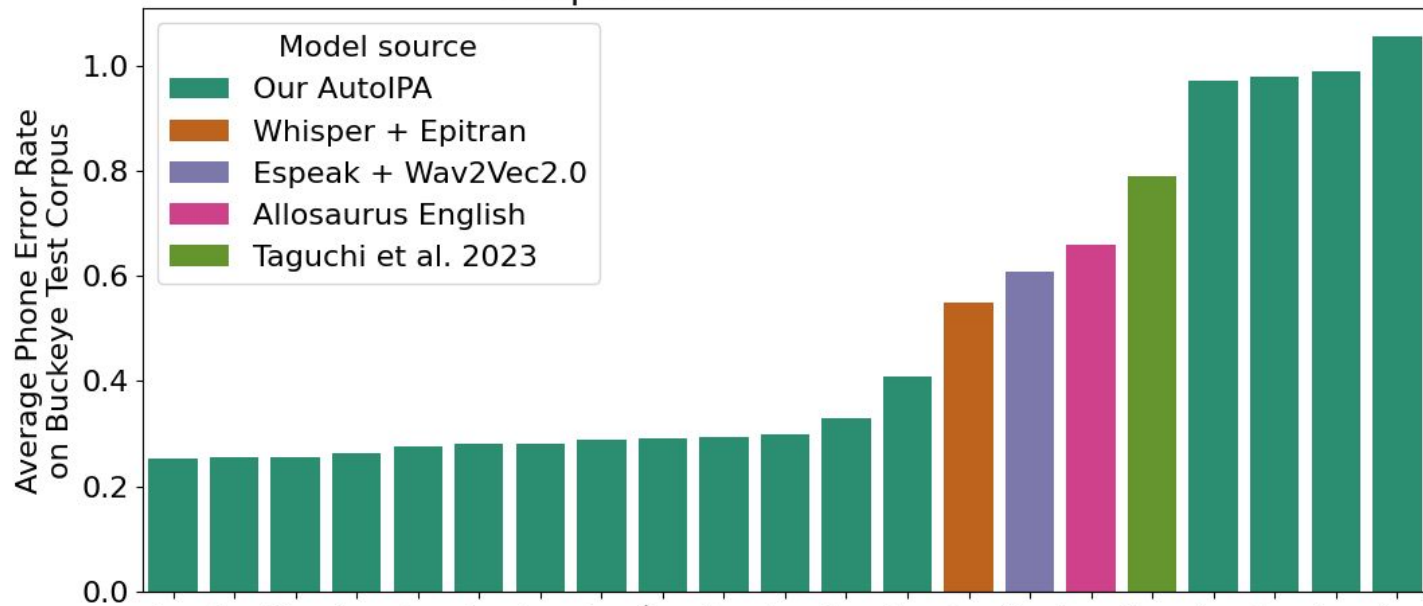
**Substitution rates, because $v \neq w$*

Relationship between Pillai score within Buckeye and substitution rate in the Buckeye test split



Impact of the quantity of fine-tuning data

Comparison of Phone Error Rates



Experiment group

Our AutoIPA: fine-tuned on 20000 samples

Our AutoIPA: fine-tuned on 12800 samples

Our AutoIPA: fine-tuned on full train split

Our AutoIPA: fine-tuned on 6400 samples

Our AutoIPA: biased female

Our AutoIPA: balanced gender

Our AutoIPA: only female speakers

Our AutoIPA: biased male

Our AutoIPA: only male speakers

Our AutoIPA: fine-tuned on 3200 samples

Our AutoIPA: fine-tuned on 1600 samples

Our AutoIPA: fine-tuned on 800 samples

Whisper + Epitran

Espeak + Wav2Vec2.0

Allosaurus English

Taguchi et al. 2023

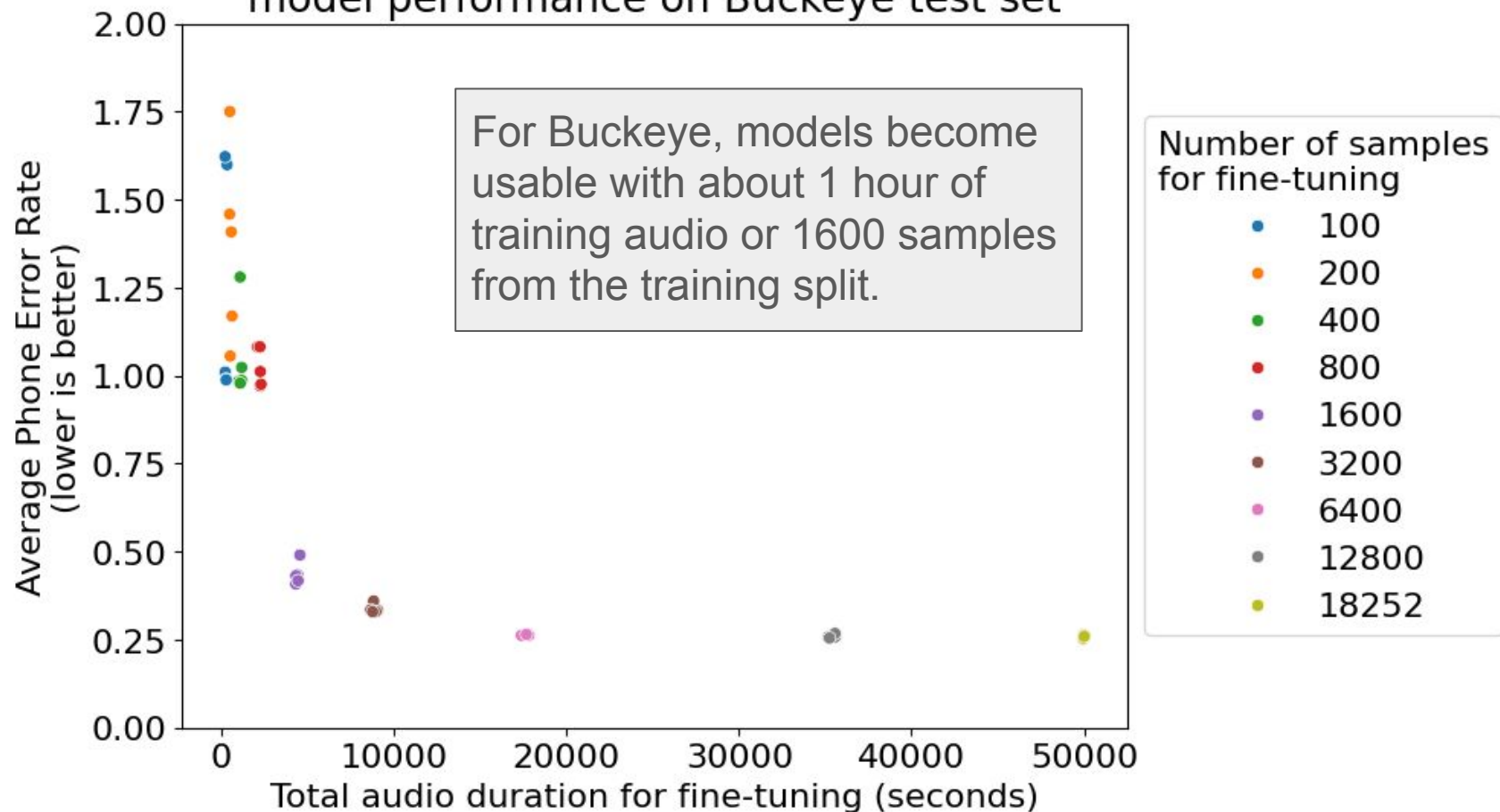
Our AutoIPA: fine-tuned on 200 samples

Our AutoIPA: fine-tuned on 400 samples

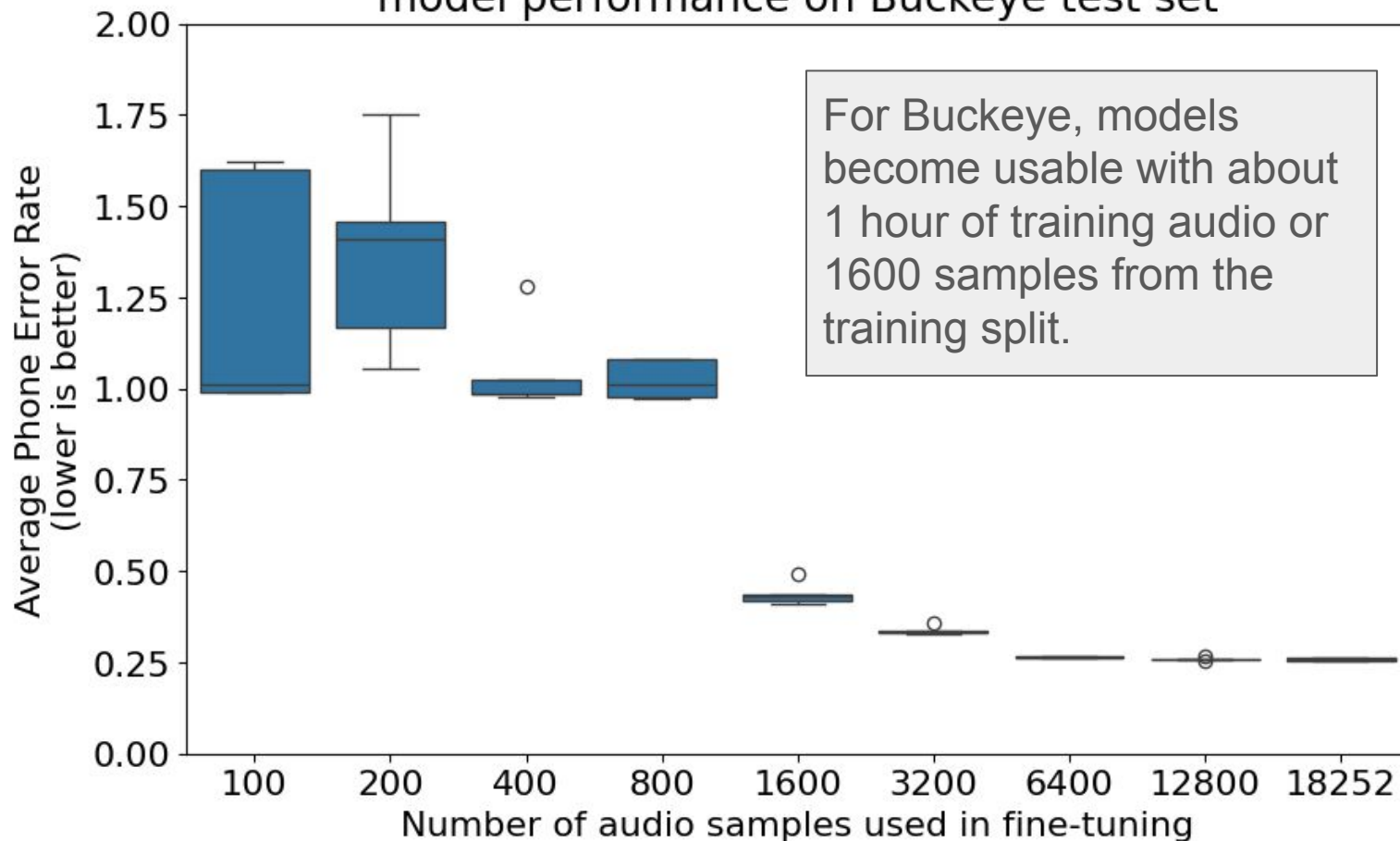
Our AutoIPA: fine-tuned on 100 samples

Our AutoIPA: fine-tuned on 200 samples

Amount of data for fine-tuning and
model performance on Buckeye test set



Amount of data for fine-tuning and
model performance on Buckeye test set



Amount of fine-tuning data: Practical implications



Speaker S25 in Buckeye test split: "over by Riverside Hospital"

Original Buckeye transcription: [oʊ v ɹ b aɪ ɹ ɪ v ɹ s aɪ d h ə s p ɪ r l]

Model fine-tuned on full Buckeye train split	I v ɹ b aɪ I ɹ ɪ v ɹ s aɪ d h ə s p ɪ r l
12800 sample fine-tuned model	ɛ v ɹ b aɪ ɹ ɪ v ɹ s aɪ d h ə s p ɪ r l
6400 sample fine-tuned model	ɛ v w ɹ b aɪ ɹ ɪ v ɹ s aɪ d h ə s p ɪ r l
3200 sample fine-tuned model	ɛ v ɹ b eɪ ɹ ɪ v ɹ s aɪ eɪ d h ə s p ɪ r oʊ
1600 sample fine-tuned model	I f █ b ɛ w ɪ v █ s aɪ d h aʊ s p ɪ r l
800 sample fine-tuned model	█ █ █ █ █ █ █ █ █ █ █ █ █ █ █ █ █ █ █
400 sample fine-tuned model	t █ █ █ █ █ █ █ █ █ █ █ █ █ █ █ █ █ █

Red = Deletion, Blue = Insertion, Yellow = Substitution

Amount of fine-tuning data: Practical implications



Speaker S38 in Buckeye test split: "oh they are they're killing on it they really are"

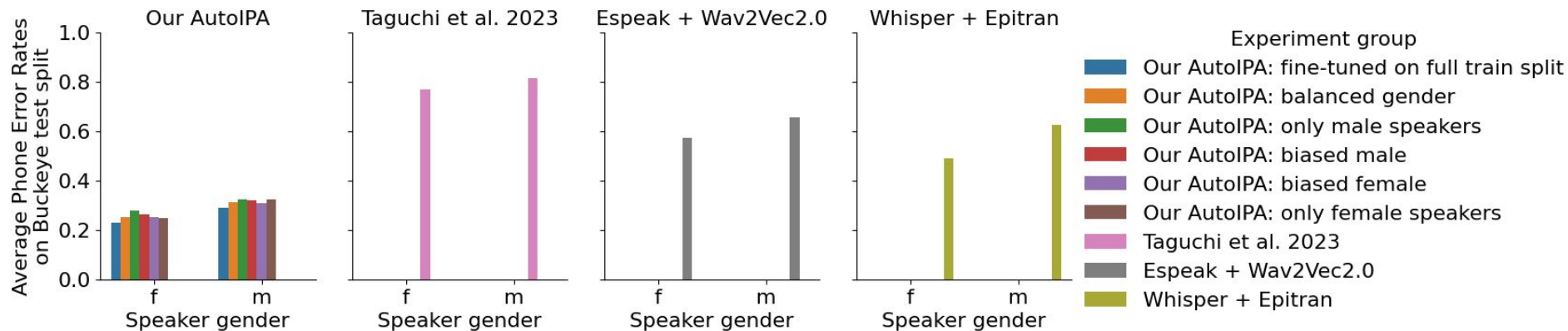
Original Buckeye transcription: [ou ð̥ɛɑɹð̥ɹkɪlɪŋɔnʌtð̥ɛɹilɪɔɹ]

Model fine-tuned on full Buckeye train split	ou ð ei a ɹ ð ε j k i l n ð ei j i l i a j s
12800 sample fine-tuned model	ou ð ei a ɹ ð ε j k i l ŋ ð ei j i l i a j s
6400 sample fine-tuned model	ou ð ei a ɹ ð ε j k i l ʌ ŋ a ð ei j i l i a j s
3200 sample fine-tuned model	ou ð ei a ɹ ð ε j k ε l ʔ n a ð ei j i l i a v d
1600 sample fine-tuned model	a ð æ a ð ε k ε l n a ʌ ε j i l a v
800 sample fine-tuned model
400 sample fine-tuned model	... t ...

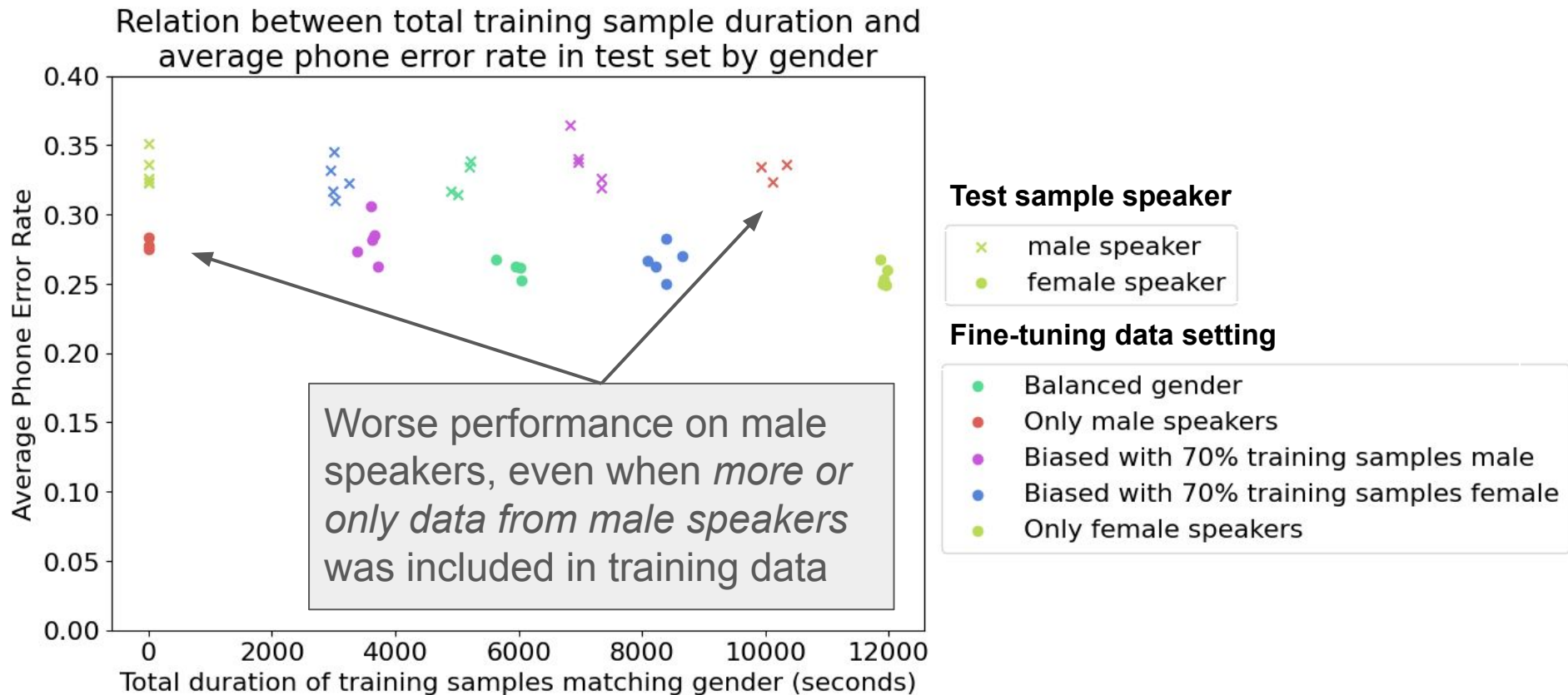
Red = Deletion, Blue = Insertion, Yellow = Substitution

Gender Effects

All models had better performance on female speakers in the Buckeye test data...



...regardless of how fine-tuning data was selected



Yao et al. 2010: 108 on Buckeye data

Generally speaking, female speakers produce slightly longer vowels than male speakers ($p=0.009$). As expected, they also have higher formant frequencies compared to male speakers ($p<0.001$ for F1, $p=0.002$ for F2). More importantly, as can be seen from Figure 3, on average female speakers have a much larger vowel space than male speakers. This is also consistent with previous findings (Byrd, 1994). Both longer duration and more expanded vowel space are indicators of clear speech (Bradlow et al., 1996), which suggests that female speakers produce clearer speech than male speakers.

Adda-Decker & Lamel 2005: on speech recognizers performing better on female speakers

Results consistently show a lower word error rate on female speech ranging from 0.7 to 7% depending on the condition. An analysis of automatically produced pronunciations in speech training corpora (totaling 4000 hours of speech) revealed that female speakers tend to stick more consistently to standard pronunciations than male speakers. Concerning speech disfluencies, male speakers show larger proportions of filled pauses and repetitions, as compared to females.

TIMIT as a test set

To begin to get a sense of how our model performs on speech from other varieties of English, we tested it on TIMIT

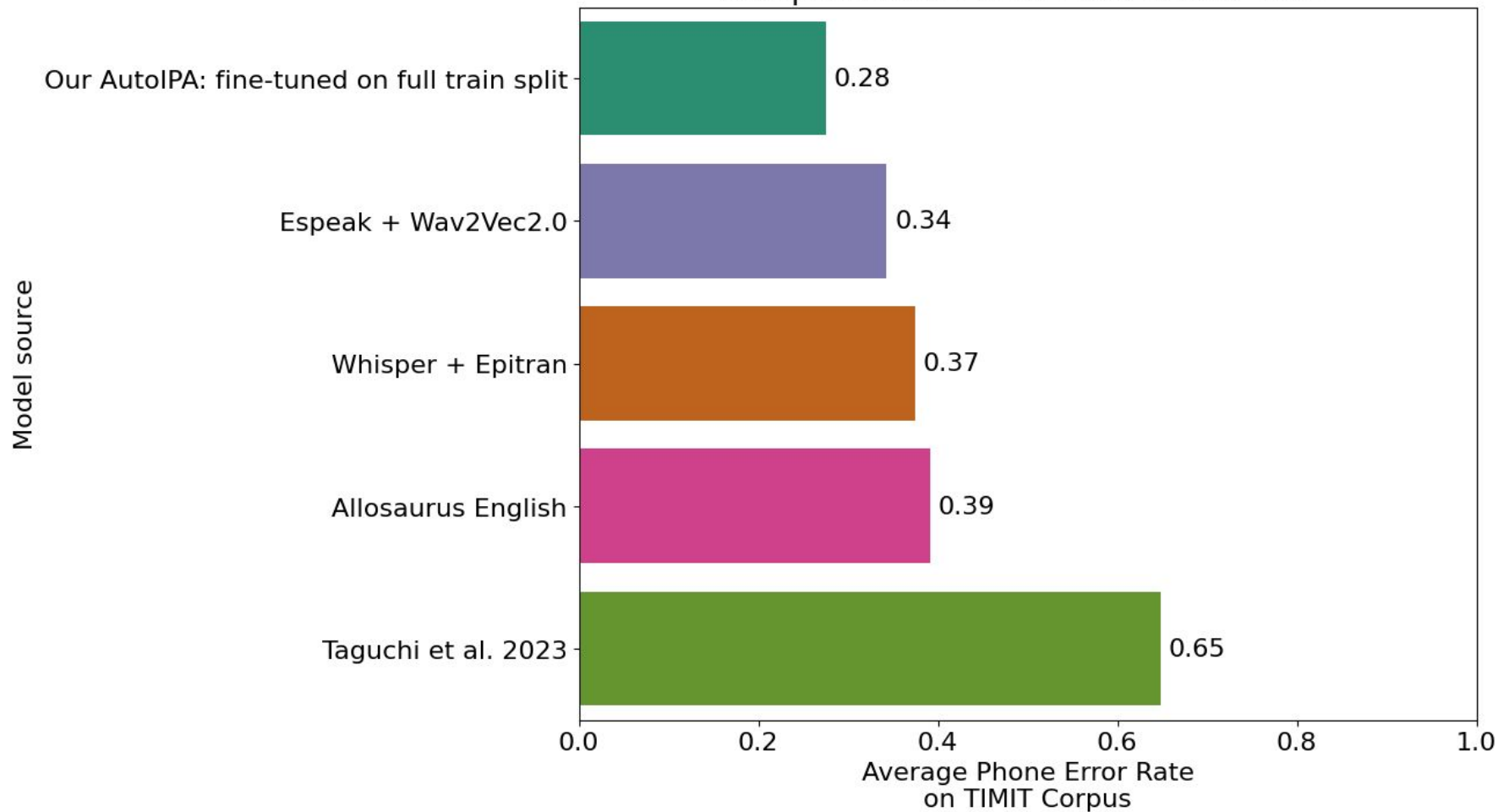
TIMIT consists of transcribed speech from 630 speakers from 8 dialect regions, each reading ten “phonetically rich” sentences. There are 5 hours of speech in total.

The ARPABET to IPA translation was trickier for TIMIT, because stop closure is indicated separately from release. We decided to merge all closure-release sequences for the same phone into a single segment.

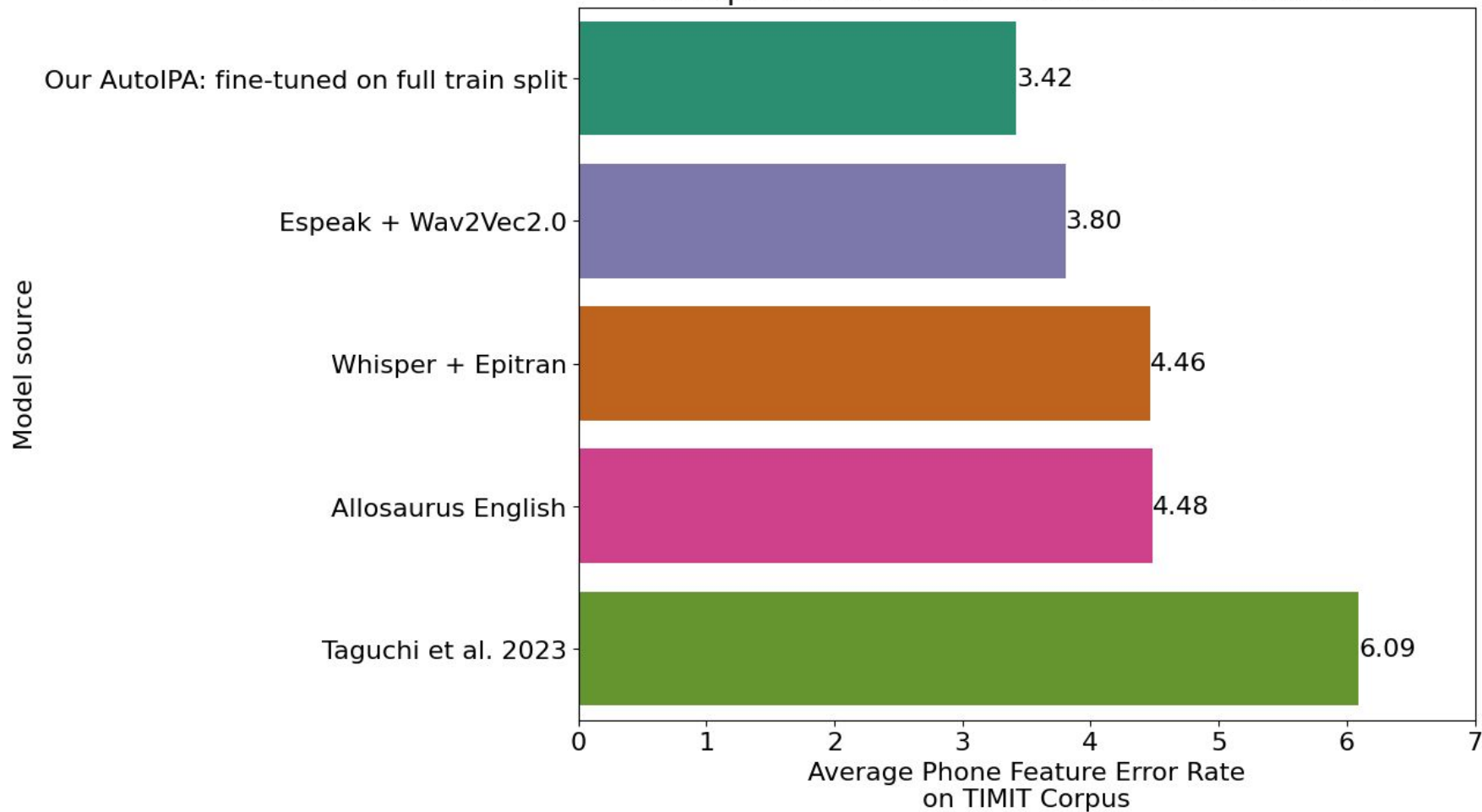
The mappings can be examined [here](#) (note that Phonocodes is greedy, so the longer input will be chosen first)

The standard approach to TIMIT in speech recognition ([Lee and Hon 1989](#)) is to merge all closures with silence. This means that all coda unreleased stops are lost (!)

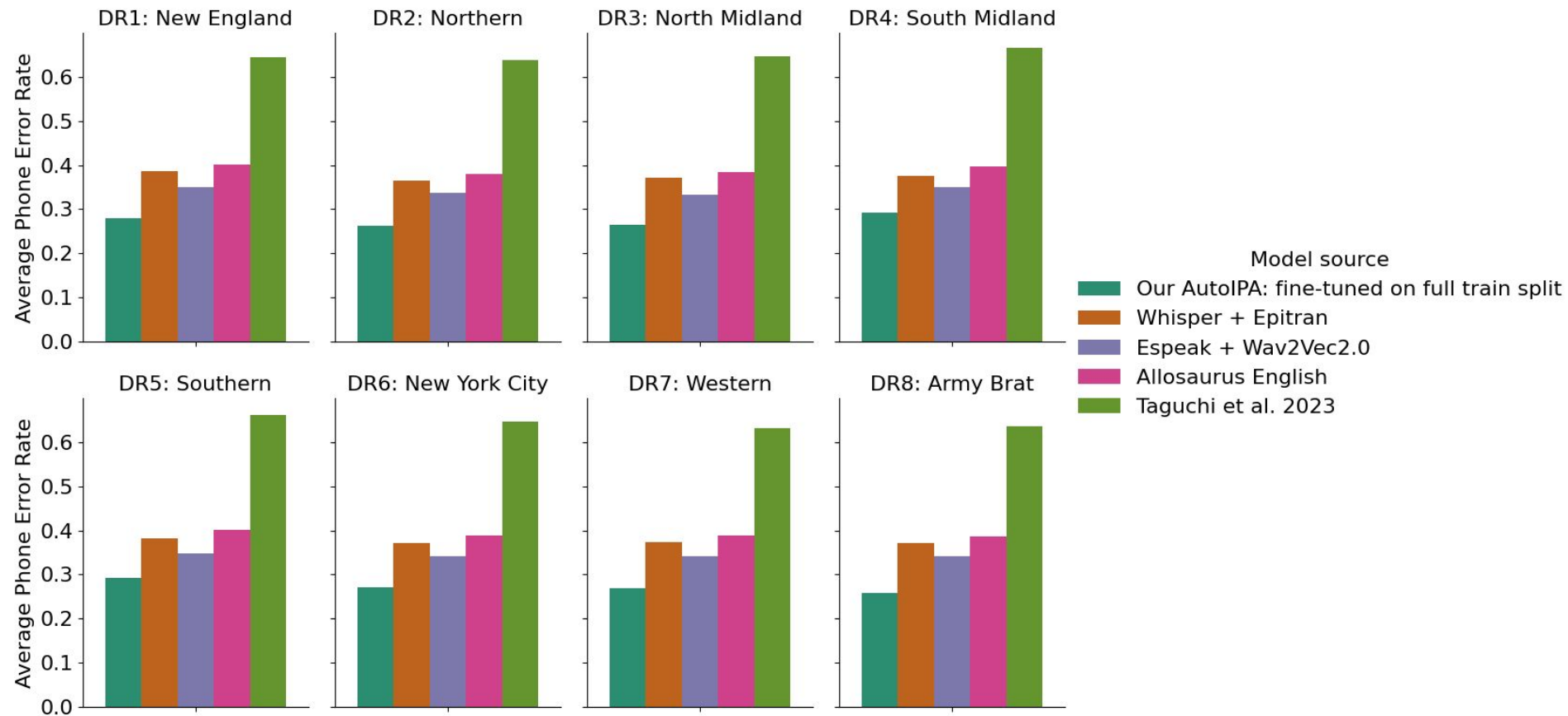
Comparison of Phone Error Rates: TIMIT



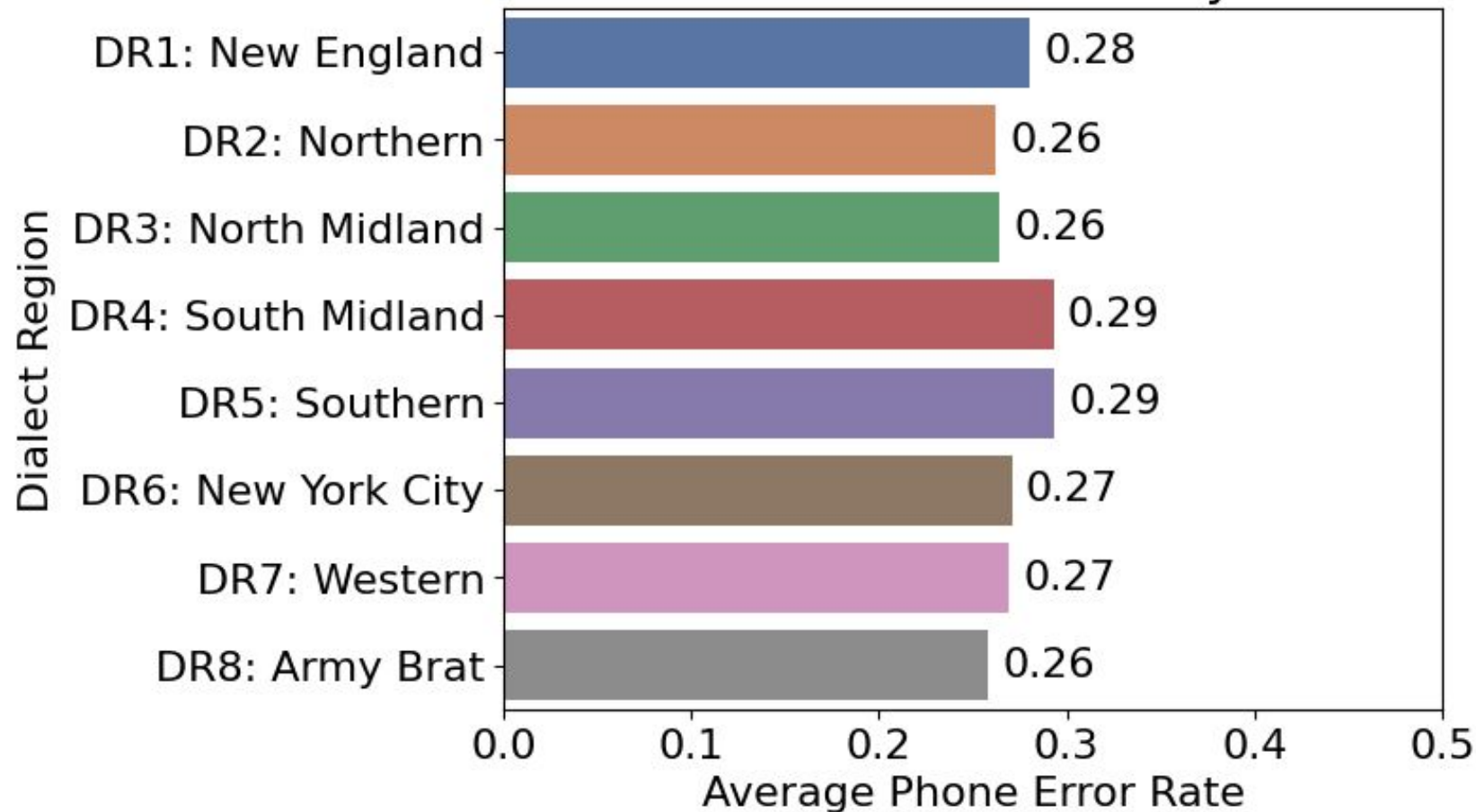
Comparison of Phone Feature Error Rates: TIMIT



Models' Average Phone Error Rates by Dialect Region



Our AutoIPA's TIMIT Performance by Dialect Region



Same sentence, different regions, TIMIT vs. AutoIPA

She had your dark suit and greasy wash water all year

New England (“best New England accent so far”) (VMH0, F, b. 1960, rec. 1986)

ʃi hæd jɪ dɑk sʌtɪŋ ɡɪsi wɑʃ wəɾə ʔɔl jɪə TIMIT

ʃi hæd jɪ dɑk suʔɪŋ ɡɪsi wɑʃ wɑɾɐ ɔʊ jɛ AutoIPA

New York (“has good NY pronunciation of ‘saw’”) (HXS0, F, b. 1941, rec. 1986)

ʃi hæd jʊ dɑk sʌtʔɪn ɡɪsi wɑʃ wɔɾɐ ʔɔl jɪə TIMIT

ʃi hæd juɐ dɑk sʌtɪn ɡɪsi wɑʃ wɔɾɐ ɔʊl jɪɐ AutoIPA

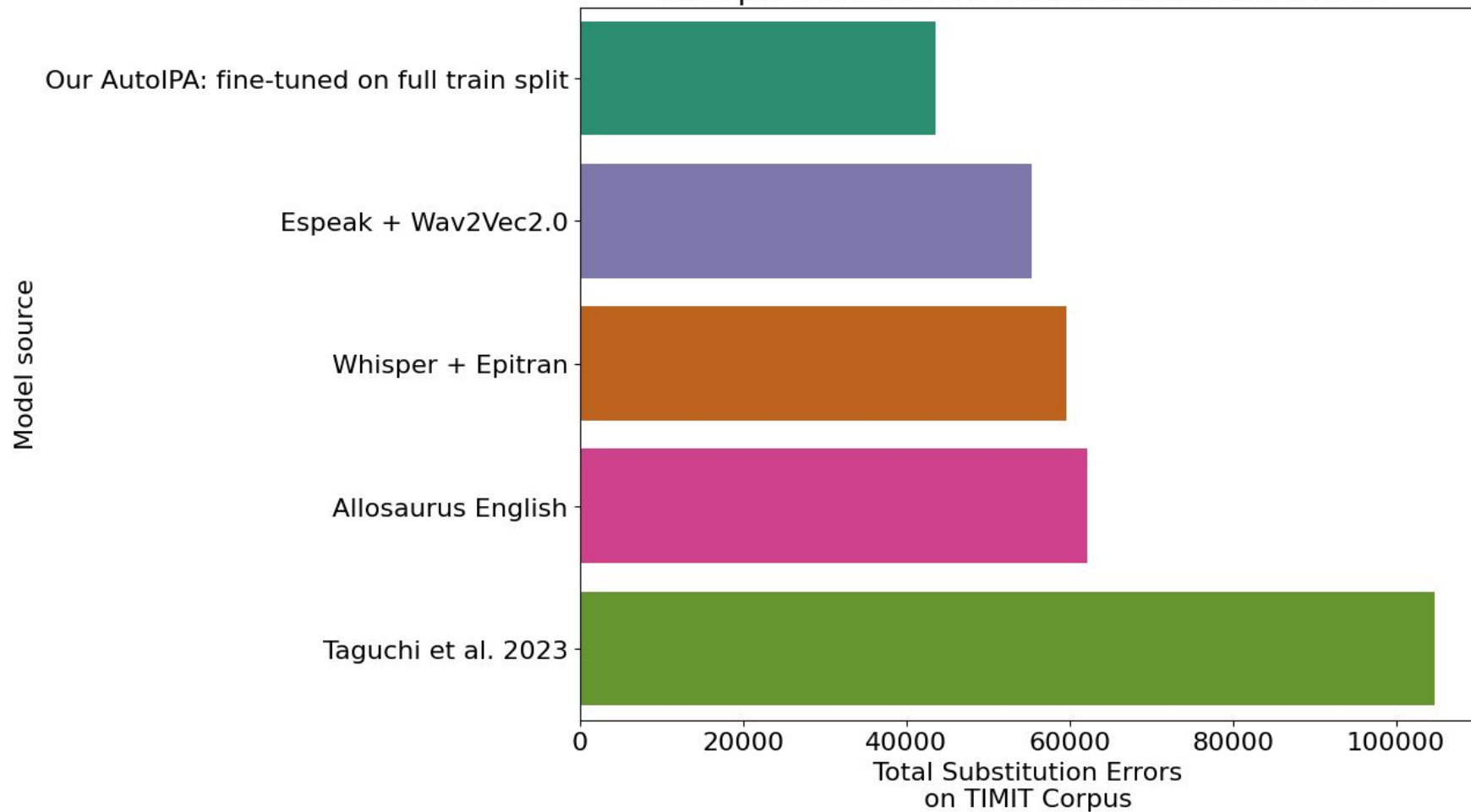
North Midland (DSS1, M, b. 1955, rec. 1986)

ʃi hæd jə dɑk sʌɾɪn ɡɪsi wɔɾʃ wɔɾə ɔl jɪə TIMIT

ʃi hæd jɪ dɑk sʌɾɛn ɡɪɛsi wɔɾʃ ʃɐɾɪ ɔl jɪɪ AutoIPA

TIMIT Error Analysis

Comparison of Total Substitutions Errors: TIMIT



Top substitution errors on the TIMIT Corpus

Our AutoIPA

Substitution	Count of error
i → I	4704
ə → ʌ	3741
æ → ɹ	3202

Espeak + Wav2Vec2.0

Substitution	Count of error
i → i:	4075
i → I	3502
ɑ → ɑ:	3435

Whisper + Epitran

Substitution	Count of error
I → j	6444
i → ə	4545
ʊ → w	3049

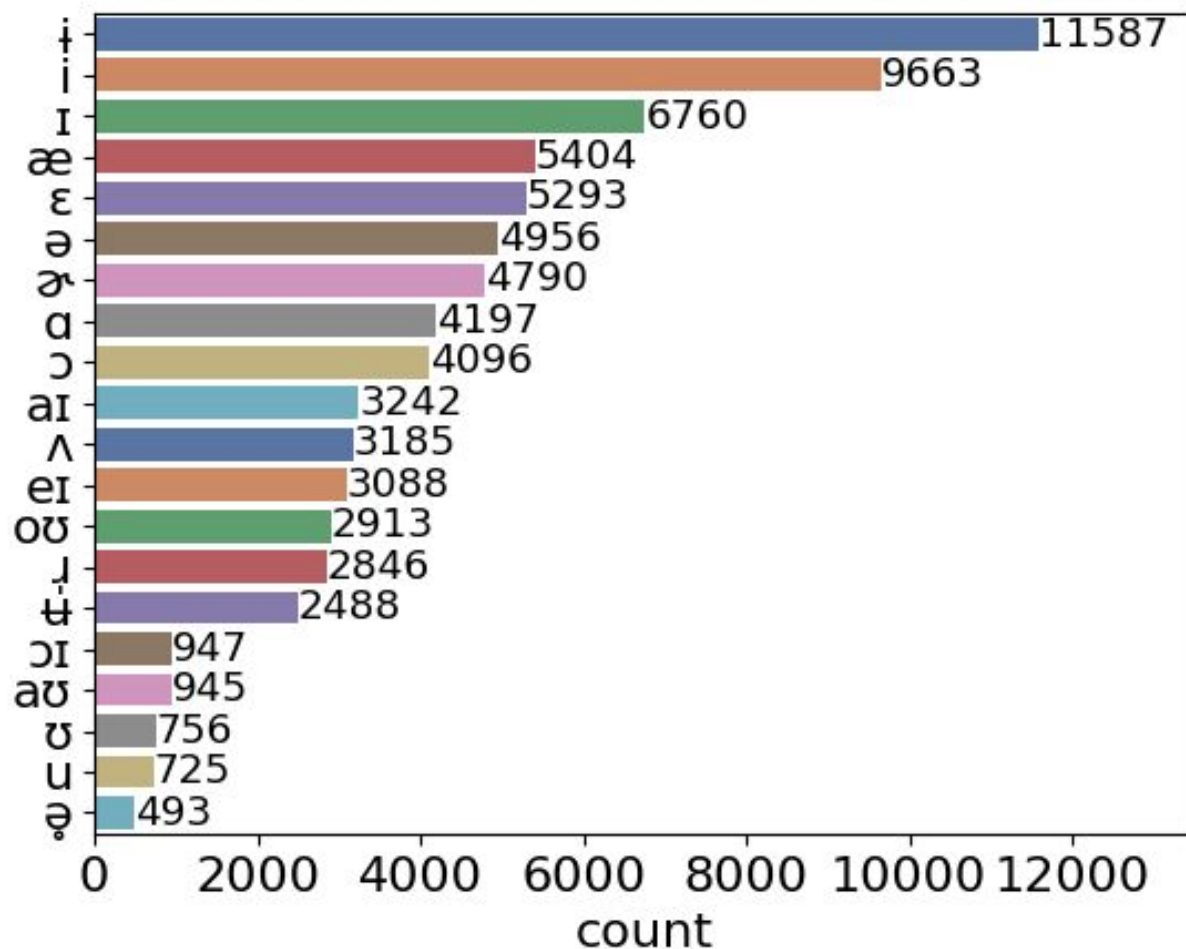
Allosaurus English

Substitution	Count of error
I → j	6189
i → ə	4515
i → I	3122

Taguchi et al. 2023

Substitution	Count of error
I → i	2777
ɹ → r	1948
I → j	1763

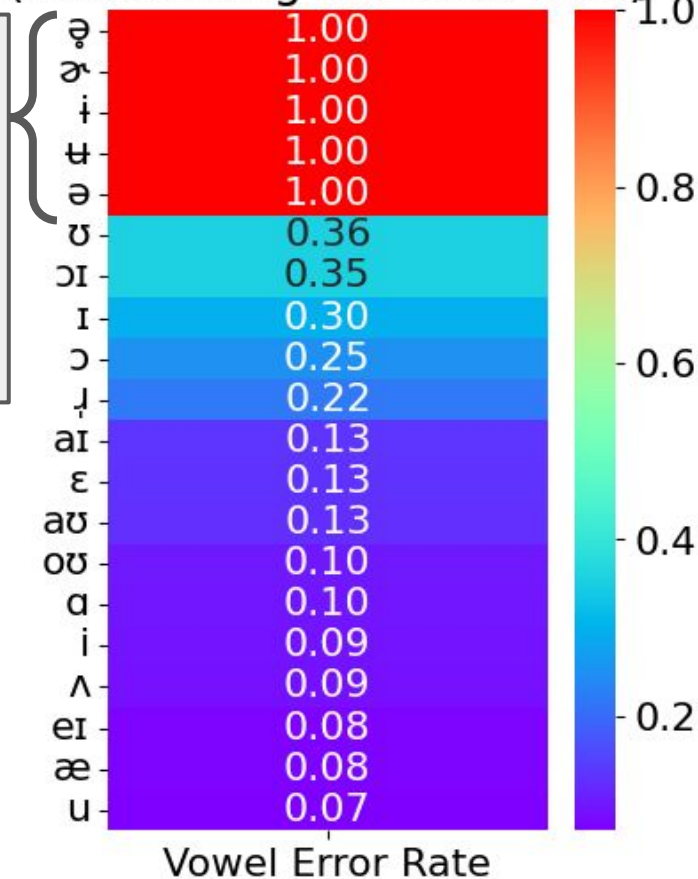
Counts of TIMIT Vowel Occurrences



These errors reflect differences in transcription conventions between corpora.

Because these are not in the Buckeye vocabulary, our model cannot output the symbols.

AutoIPA TIMIT Vowel Error Rates
(Descending worst to best)



Distributions of Errors for a Given Vowel

For a given vowel v in the TIMIT vocabulary V , we want to understand **how the model is getting the vowel wrong**.

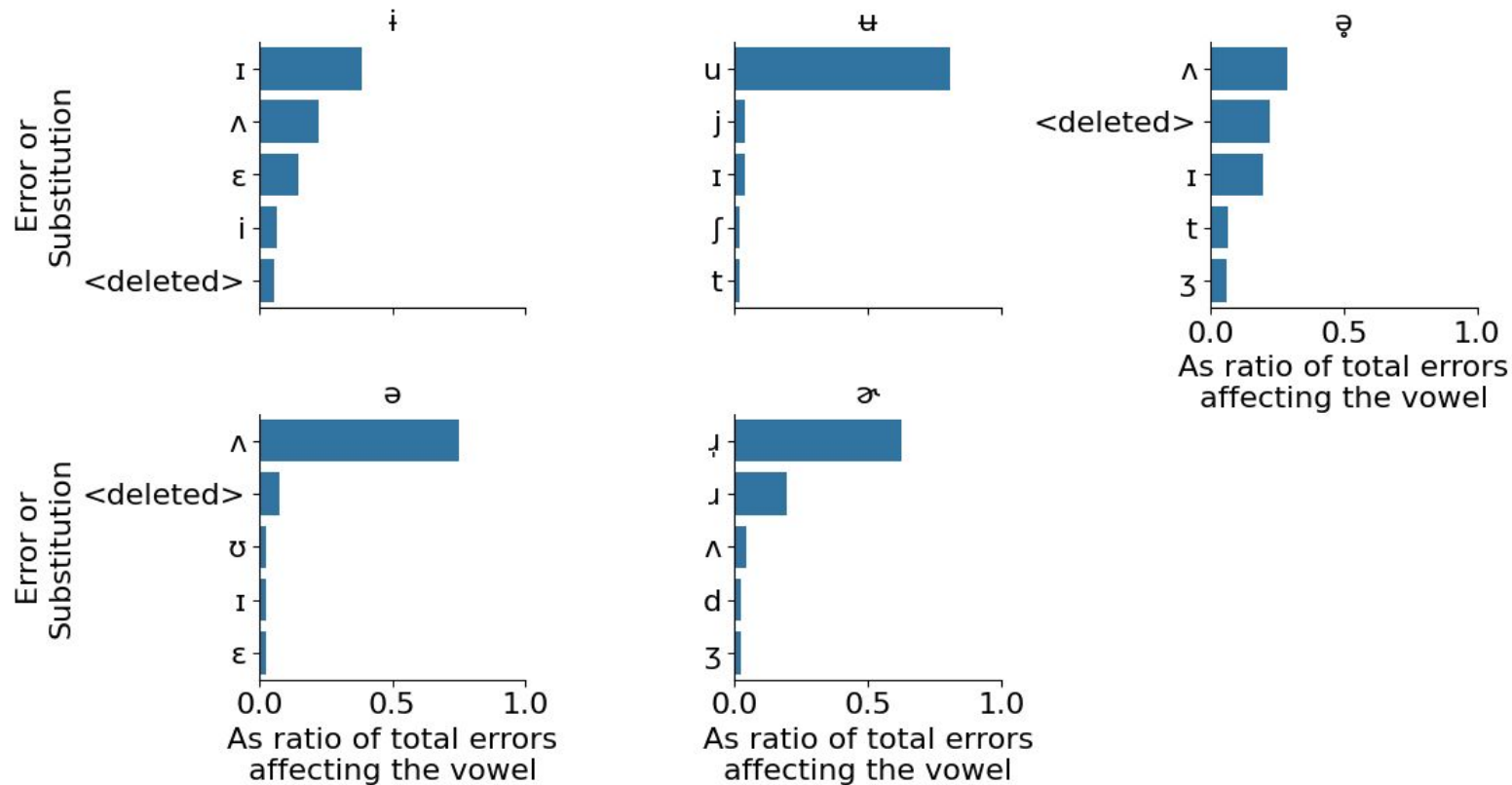
What **kinds of mistakes** and **in what distributions** does the model make for each vowel? Can this tell us anything about the corpus or its language varieties?

For a fixed vowel v and specific type error on v (deletion or substitution by w), normalize by the total number of errors:

$$\frac{\text{count_deletion}(v)}{\text{count_errors}(v)} \qquad \frac{\text{count_substitutions}(v,w)}{\text{count_errors}(v)}$$

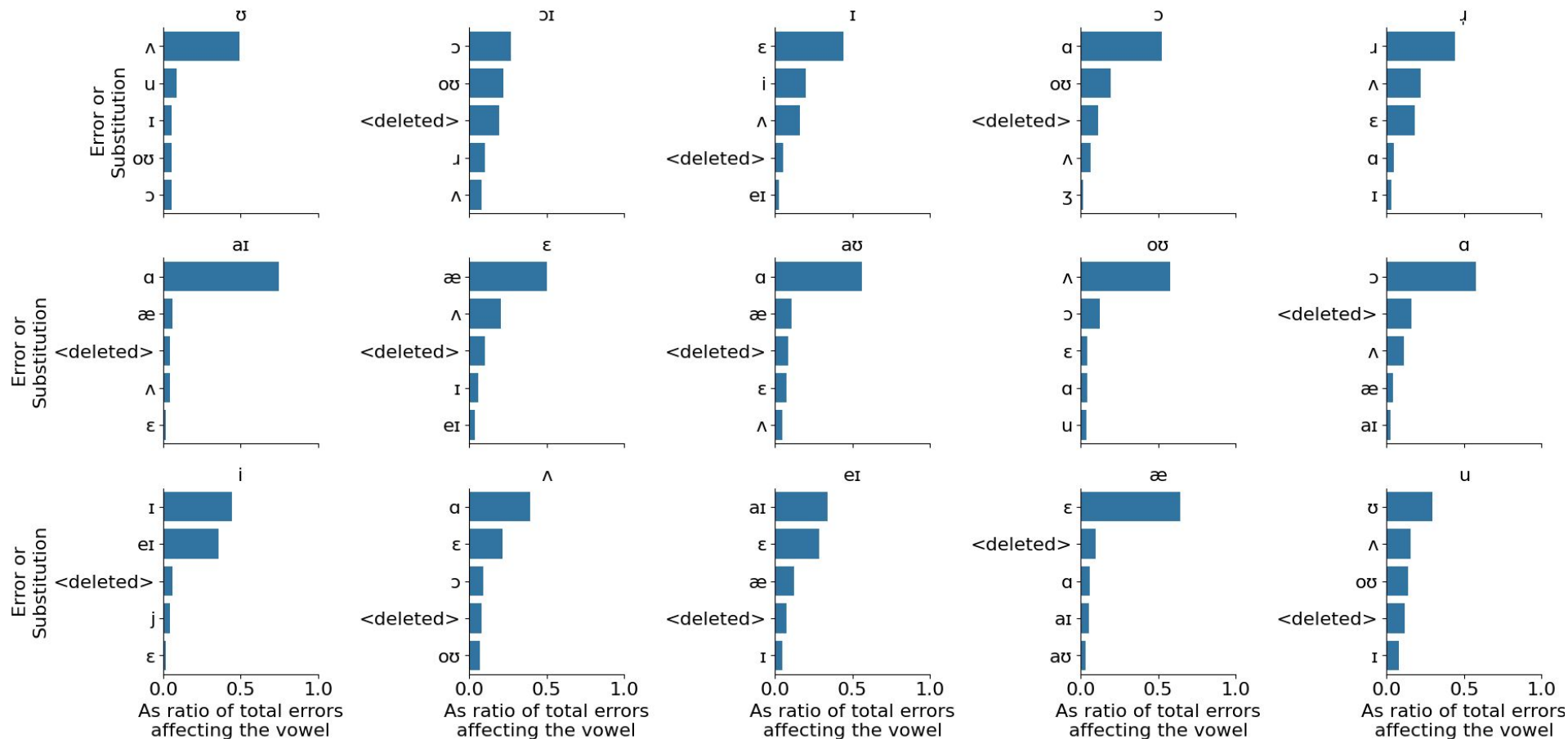
** Conditional probability of each error given that there is some error affecting v*

Top 5 errors for vowels AutoIPA always incorrectly transcribes

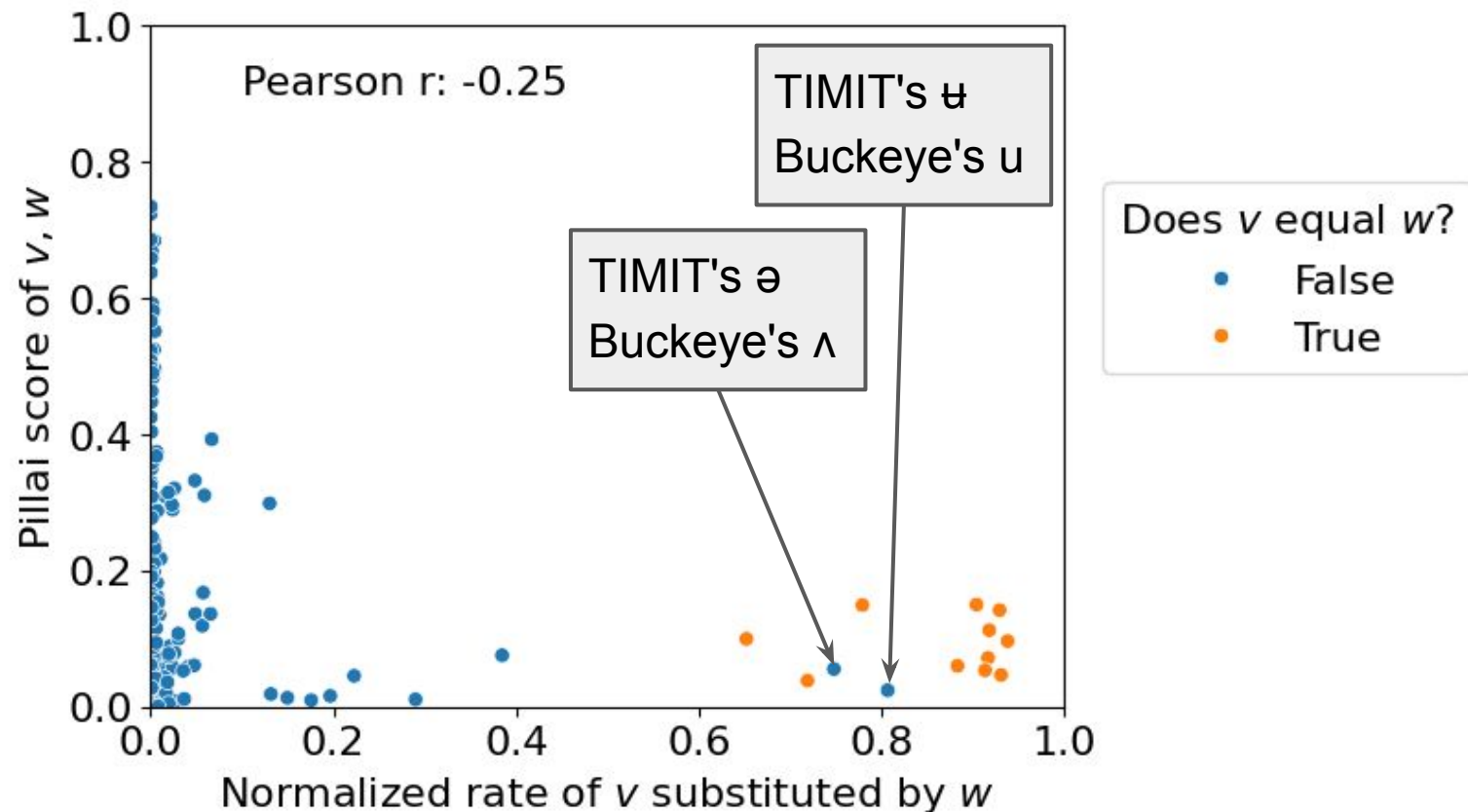


Correctable with simple replacement in pre- or post- processing

Remaining TIMIT Vowels: Top 5 AutoIPA Errors for each vowel

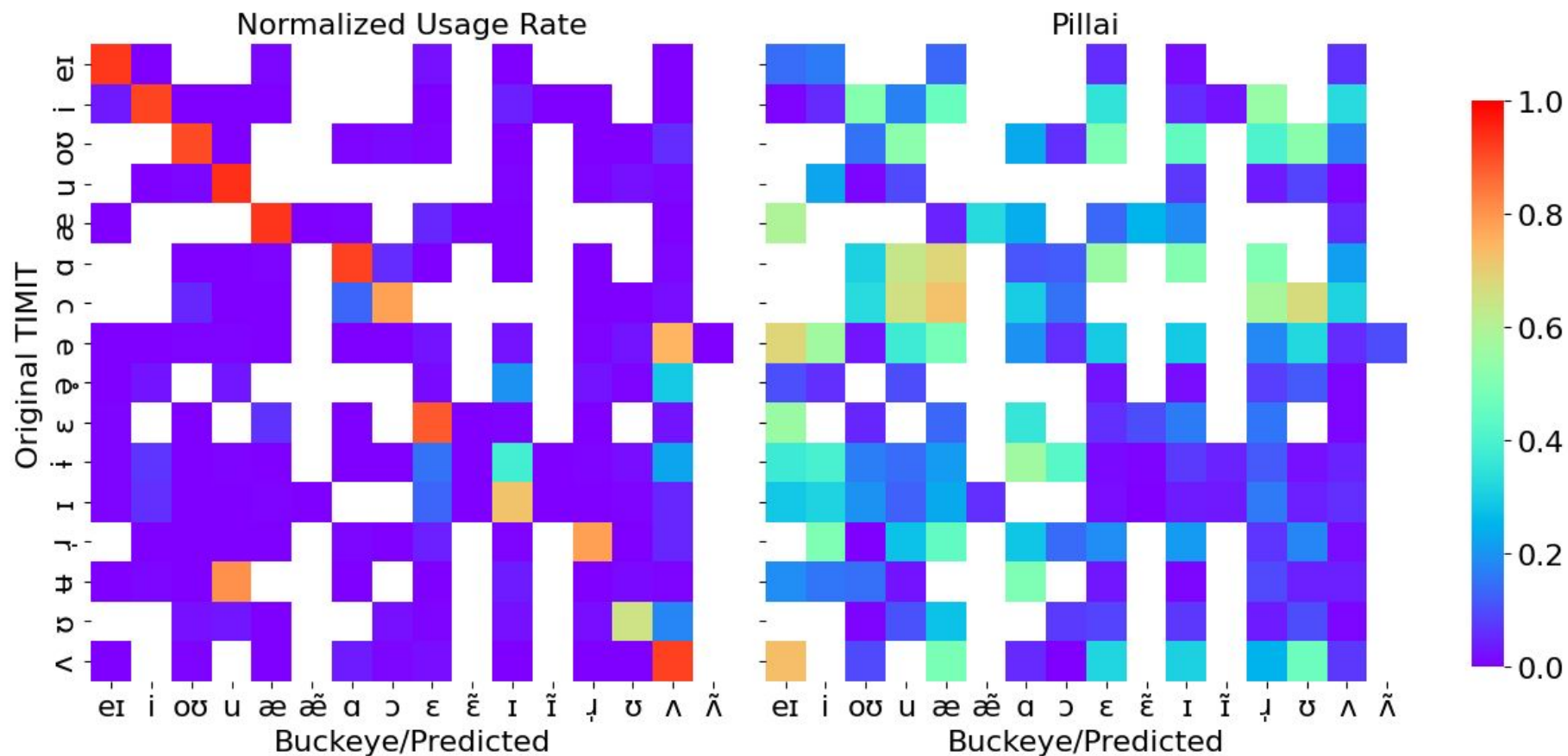


Pillai scores of TIMIT vs. Buckeye Vowels and normalized usage rates by AutoIPA on TIMIT



**Usage rates, because it's possible that $v=w$*

Pillai scores of TIMIT vs. Buckeye Vowels and normalized usage rates by AutoIPA on TIMIT



Web-based implementation with text grid support

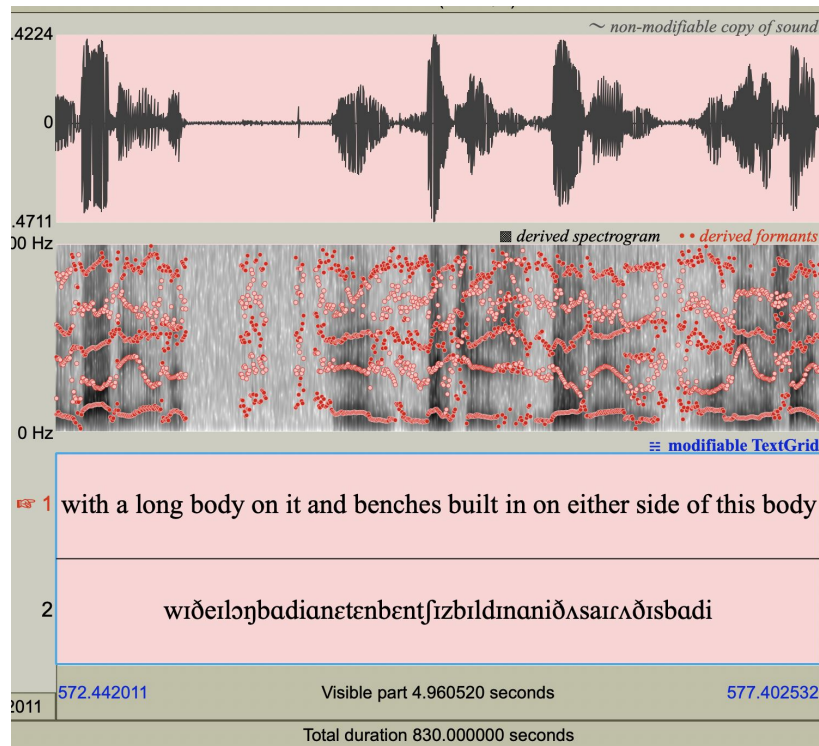
AutoIPA is now [available on hugging face](#), with support for Praat text grids

We are now working on outputting phone alignments

Transcription of JP (quasi-Canadian) saying
“I put my **cat** on a **cot**. I put my **cat** on a **cot**.
I put my **cot** on a **cat**. I put my **cot** on a **cat**”:

aɪ pʊʔ maɪ kɑ̃ an ɫ kɔ̃ ʌ pʊʔ m kɑ̃ an nɫ kɔ̃ʔ
ʌ pʊʔ maɪ kɔ̃t an kɑ̃ʔ ʌ pʊʔ maɪ kɔ̃ʔ an nɫ kɑ̃p

It uses [æ] for JP’s “ran”.



Next steps

Wav2Vec 2.0 has been used as a classifier by [Kim et al. 2024](#) for nasalization and by [Tanner et al. 2025](#) for stop realizations

AutoIPA will likely be useful as a pretrained model to be fine-tuned for that kind of work

It likely also has many applications out of the box (especially in the study of word-final consonant realizations)

Our next step is in applying it to the study of phonological variation at the “border” of western and eastern New England dialect regions in Western Massachusetts

We foresee an iterative approach, getting first pass transcriptions from AutoIPA, correcting them, and then fine-tuning our model with the new transcriptions

[Notes](#) on the standard TIMIT phone reduction protocol

Acknowledgments

We gratefully acknowledge the support of the Center for Data Science and Artificial Intelligence, the HFA/CICS Collaborative Seed Fund, and the Public Interest Technology Initiative, all at UMass Amherst.