

Modeling Non-Native Pronunciation

Alexander Metzger | Aruna Srivastava


PAUL G. ALLEN SCHOOL
OF COMPUTER SCIENCE & ENGINEERING

48% of all foreign speakers struggle with their accent.

Goal: use phonetics to provide foreign speakers granular pronunciation feedback

Issue: most transcriptions are highly inaccurate for non-standard speech

native speech 🇺🇸

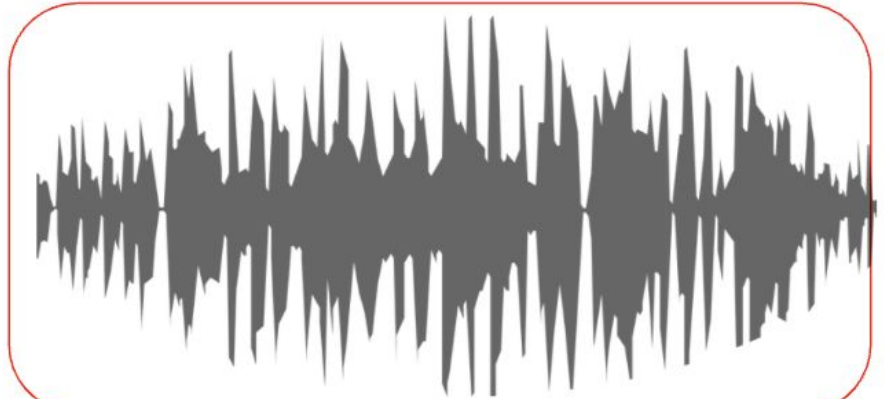


Ground truth: **kalɪŋ kaɪdz**

kəlɪŋ kuɪdz

Error: 10%

Non-native speech 🇩🇪



Ground truth: **kəlɪ karts**

kəlɪŋ kadz

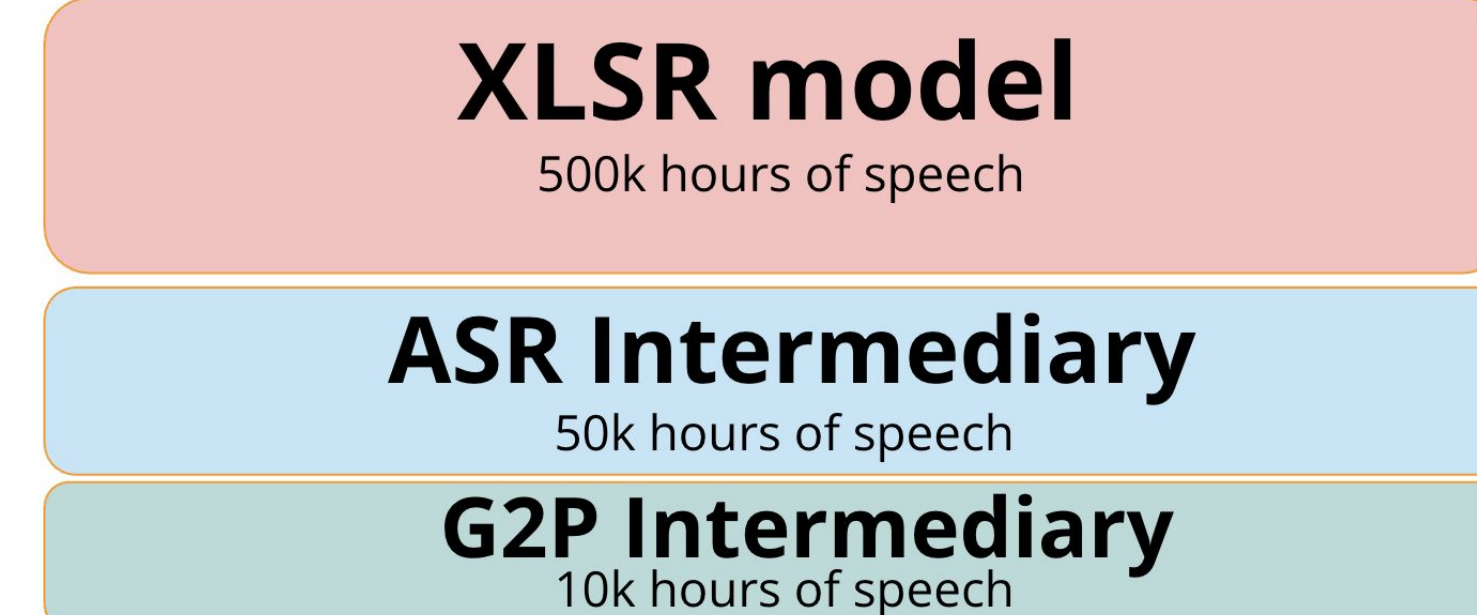
Error: 45%

Question

How can we improve a pretrained model that overfits to representing standard speech to transcribe non-standard speech phonemes, including L2 accented speech and speech impediments?

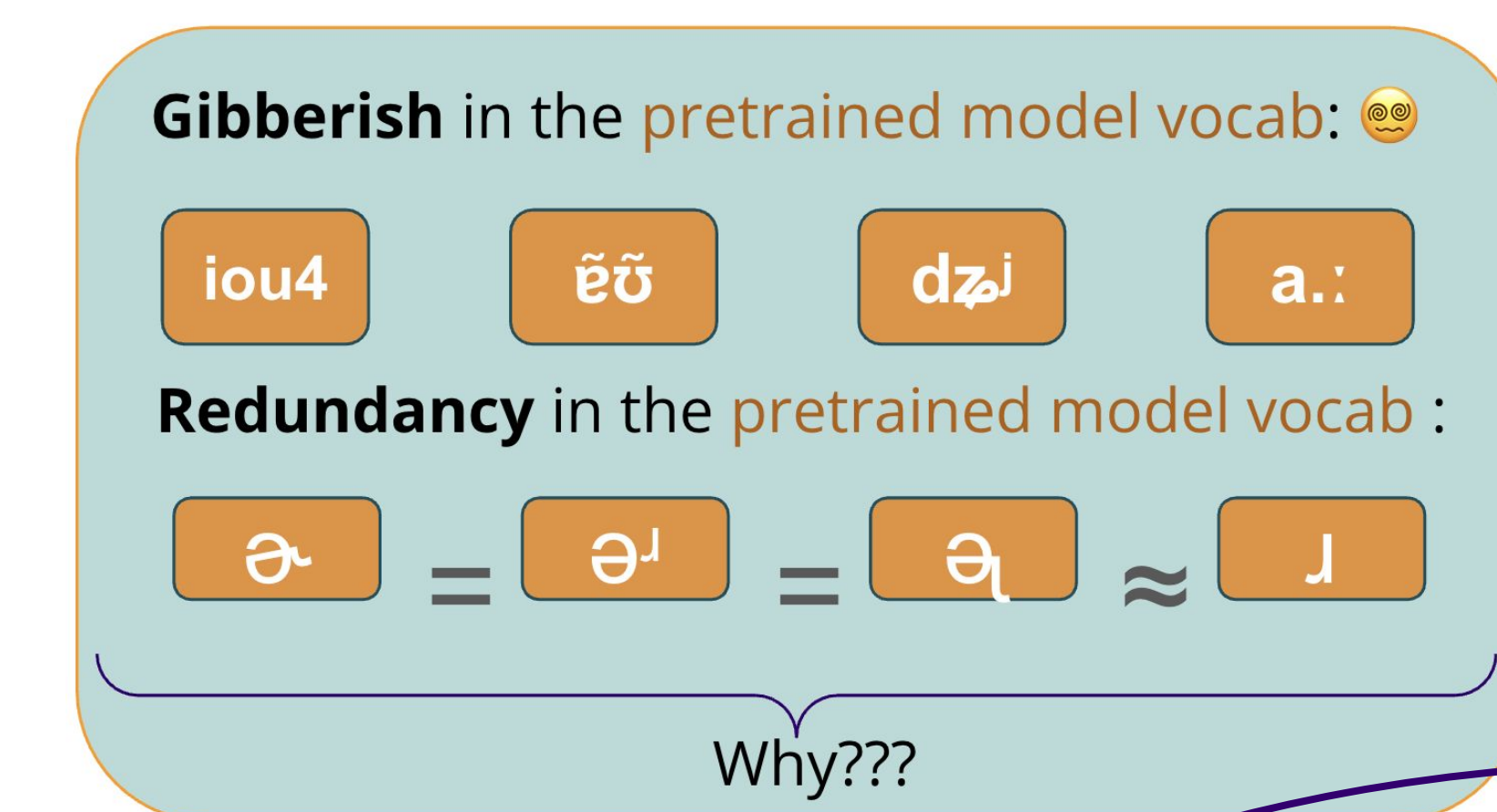
- 1) What is the **smallest set of speech sounds** needed to capture both **standard and non-standard** American English?
- 2) Does following a **curriculum** approach using computer-labeled data **prior** to limited human-labeled data improve performance?

Prior Pretrained Model



G2P Model = BAD

G2P trained checkpoint is ineffective to model non-standard speech well



Data

We curate **~100 hours** of speech data across **9** representative annotated datasets

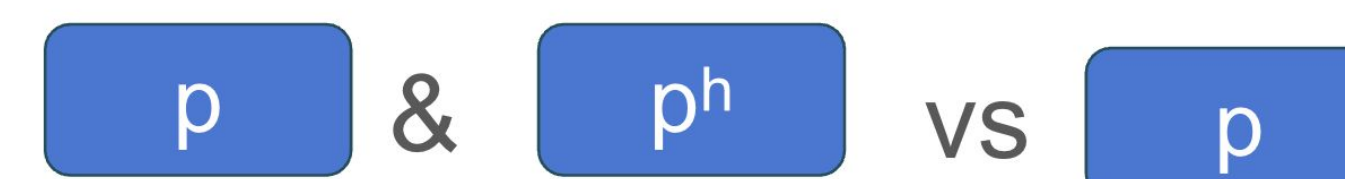
Dataset Vocab Refinement

Phonetic annotation **varies** across **every dataset and linguistic** annotator.

Datasets use **different phonetic alphabets**



Datasets use **annotators with different training**



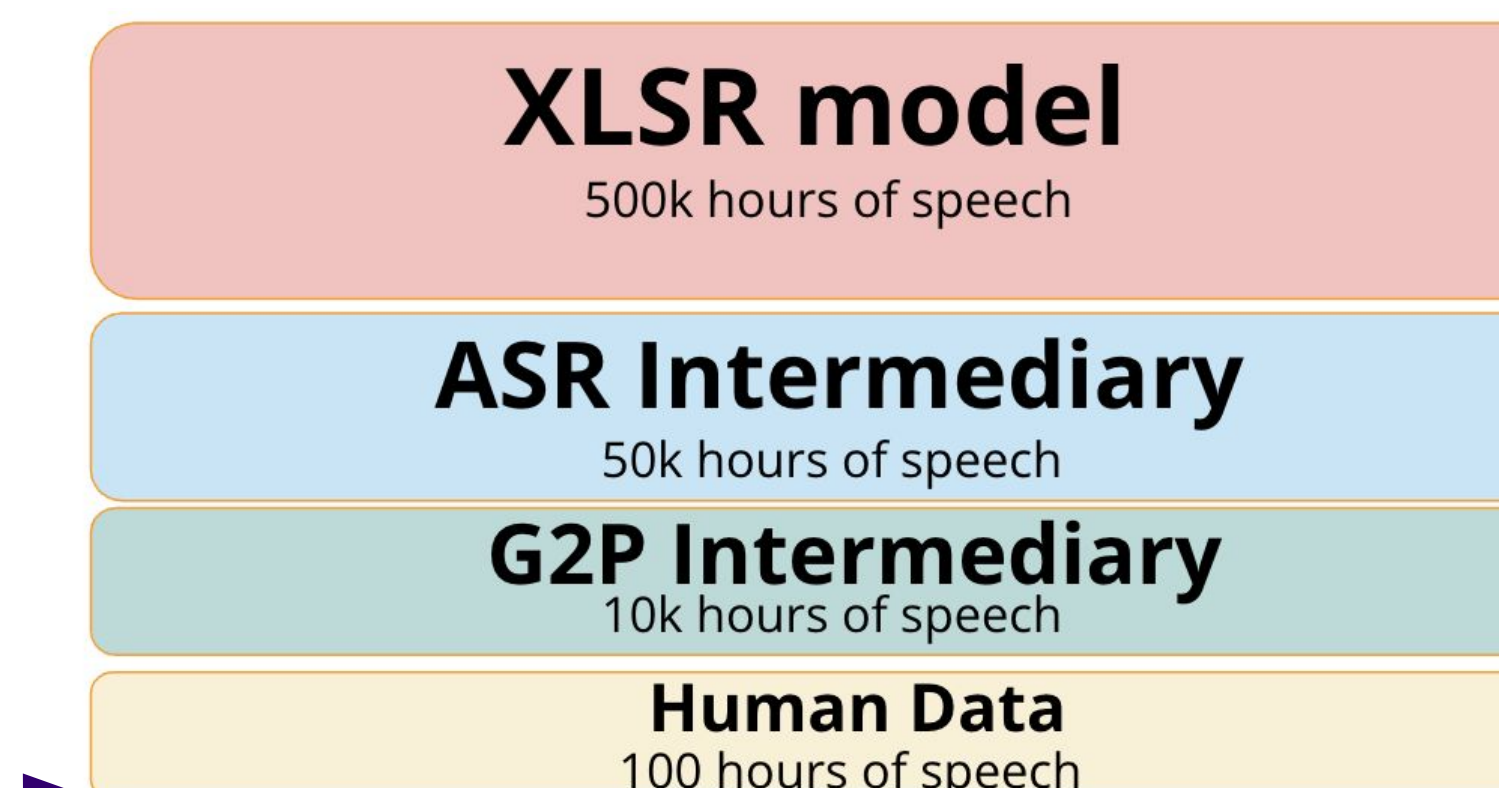
Datasets use **niche symbols**



We must standardize to IPA, collapse redundant phones, and manually remove highly bias/ambiguous samples

Our Method

Use high-quality **human annotated data** as the **final layer** of learning for phonetic transcription



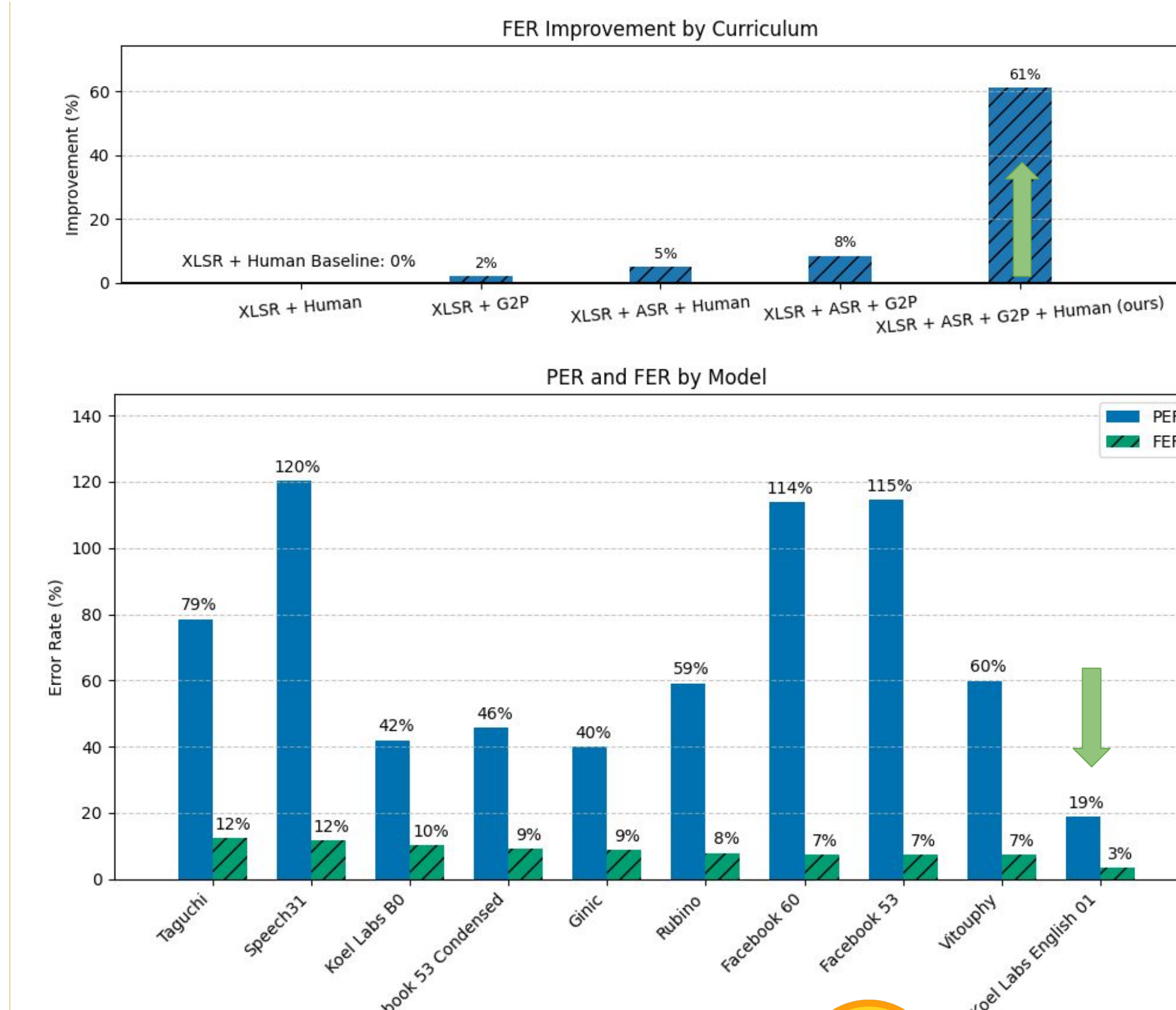
Evaluation

We use phoneme error rate (**PER**) and weighted phoneme error rate (**FER**) for evaluation.

PER considers all character differences equally, FER considers the differences by linguistic distance. For a model that has a similar vocabulary to the test set, it's PER will be superficially high compared to a model that has a slightly different vocab/phoneme notation.

PER	Weighted PER
POP	POP
BOP -1	BOP -0.25
IOP -1	IOP -0.95

Results



Open Sourced!

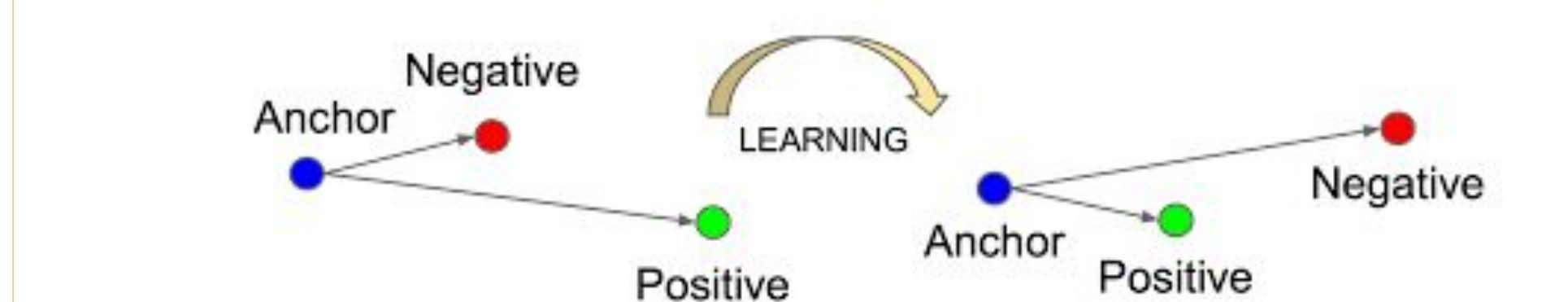
You can test our State-of-The-Art model yourself! <https://shorturl.at/f8Y1E>

Limitations

- Wav2Vec2 architecture constraints
- Speaker diversity bias
- Unaddressed annotator bias

Future Work

- Handle annotator bias
- Explore IPA representations
- Employ contrastive learning!




Interested in learning more? Sign up for the Beta!



sounds in the pretrained model

W

native speech 🇺🇸

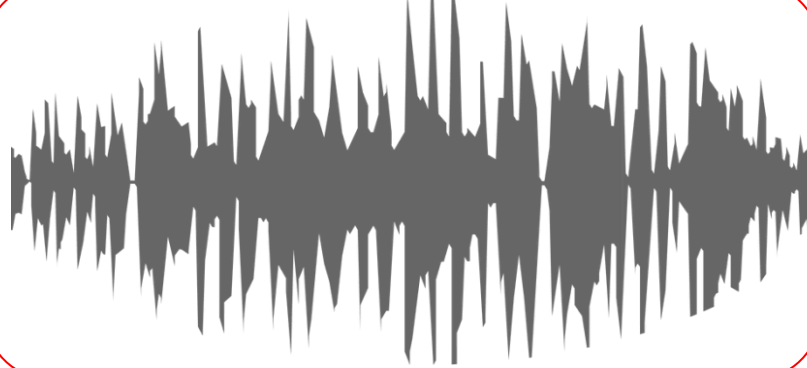


Ground truth: **kəlɪŋ kɑɪdz**

kəlɪŋ kuɪdz

Error: 10%

Non-native speech 🇩🇪



Ground truth: **kəlɪ karts**

kolɪŋ kadz

Error: 45%

ə̃ = ə' = ə̣ ≈ ɹ

Gibberish in the pretrained model vocab: 😊

iou4

ēū

dʒɪ

aː

Redundancy in the pretrained model vocab :

ə̃

 =

ə'

 =

ə̣

 ≈

ɹ

Why???

Datasets use different phonetic alphabets

@`

 =

ə̃

 =

AXR

Datasets use annotators with different training

p

 &

pʰ

 vs

p

Datasets use niche symbols

ɹ̥

gʏ

β

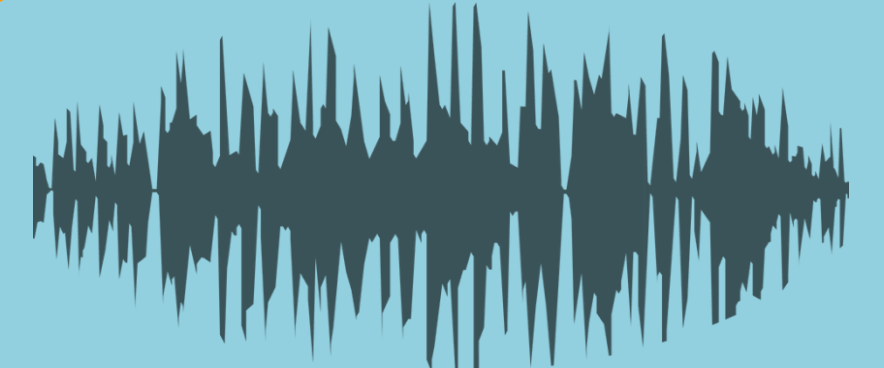
XLSR model
500k hours of speech

ASR Intermediary
50k hours of speech

G2P Intermediary
10k hours of speech

Human Data
100 hours of speech

G2P Intermediary



Transcription: **Calling Cards**

Grapheme to Phone: **kəlɪ karts**