

Self-supervised learning for speech processing

Facebook AI Research



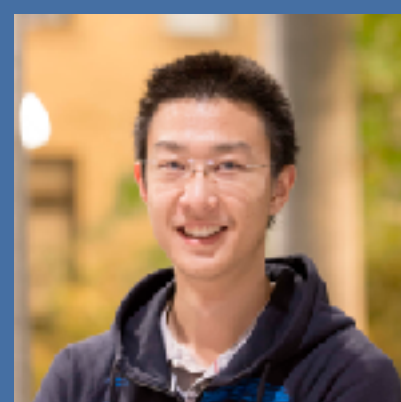
Alexei Baevski



Alexis Conneau



Steffen Schneider



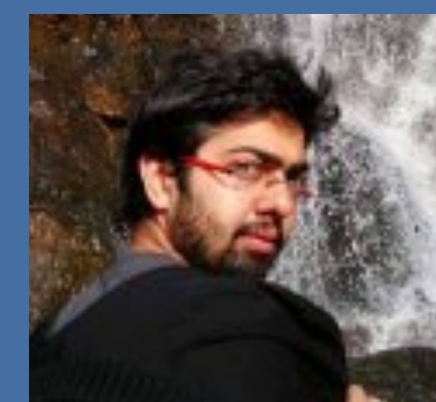
Henry Zhou



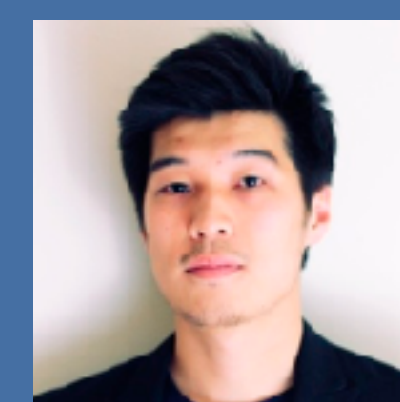
Abdelrahman
Mohamed



Anuroop
Sriram



Naman
Goyal



Wei-Ning Hsu



Michael Auli



Kritika Singh



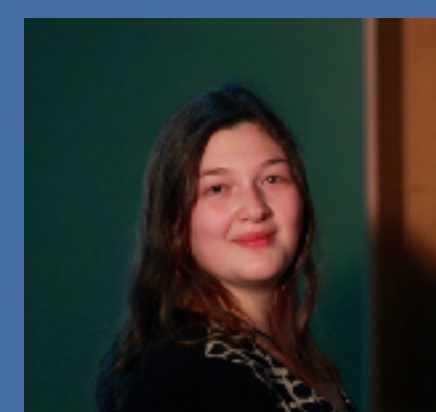
Yatharth Saraf



Geoffrey Zweig



Qiantong Xu



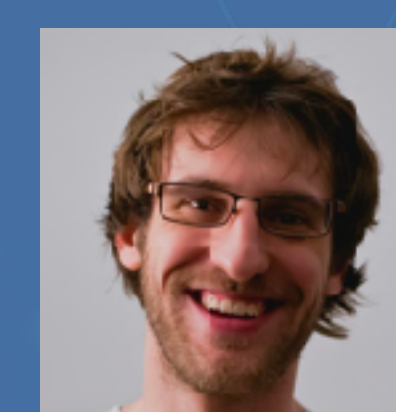
Tatiana
Likhomanenko



Paden
Tomasello



Ronan
Collobert



Gabriel
Synnaeve

Speech technology



Video captioning

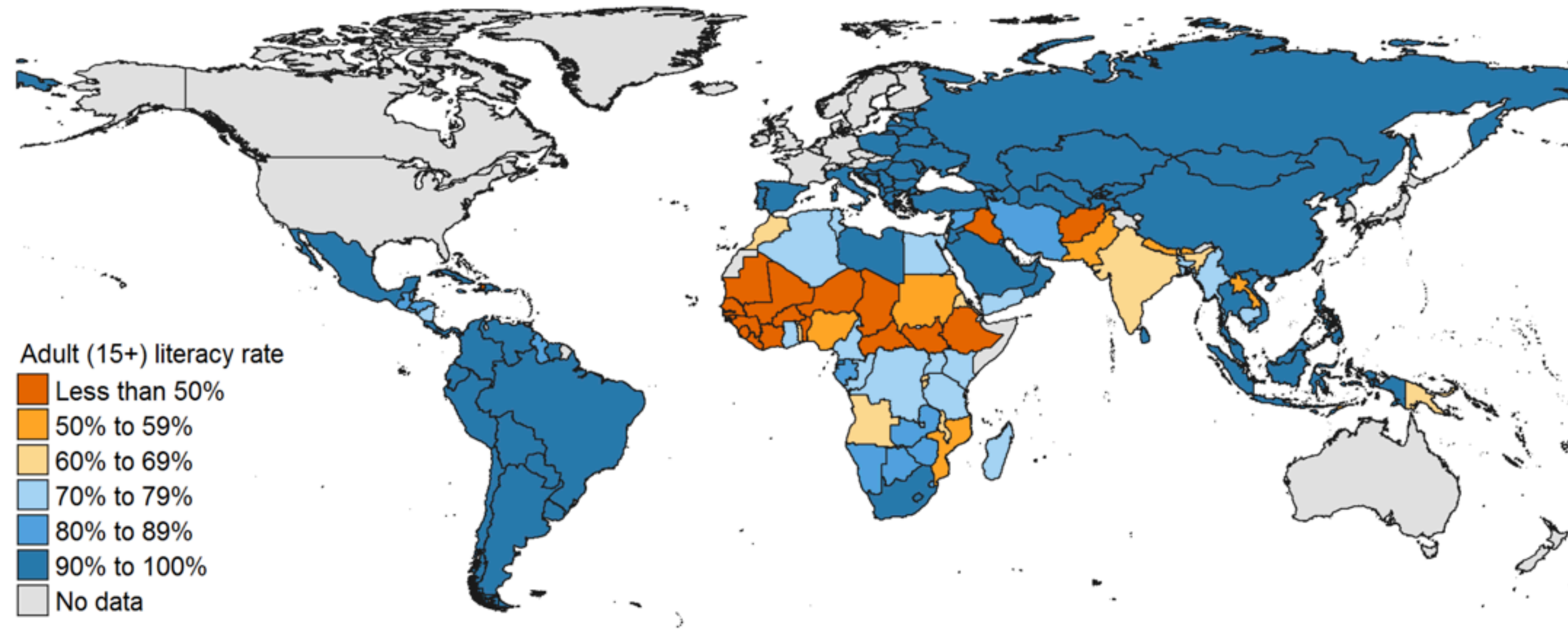


Mobile devices



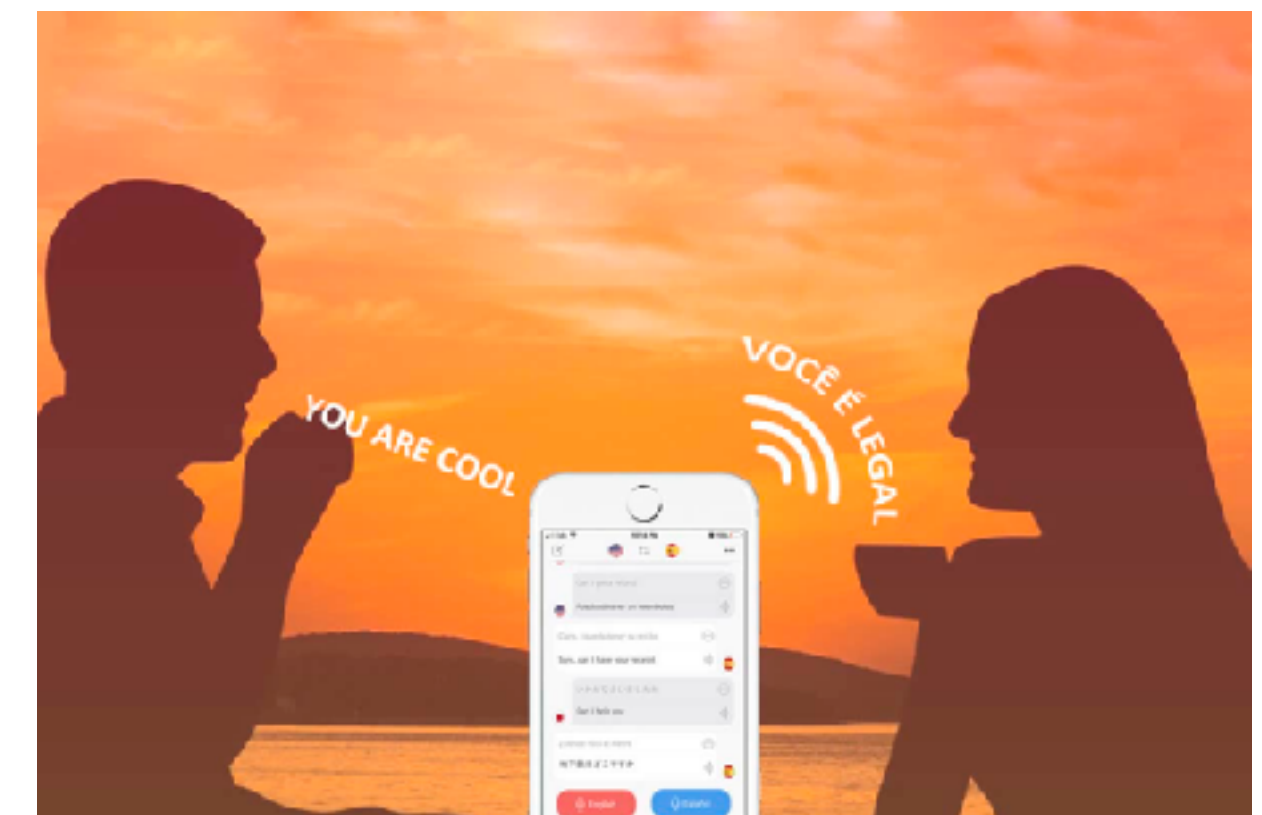
Home devices

Adult literacy rate by country, 2016



Speech applications

- Speech to text/speech recognition - dictation etc.
- Text to speech - reading out aloud
- Keyword spotting - "Hey Alexa/Portal"
- Speaker identification - is it your voice?
- Language identification
- Speech translation



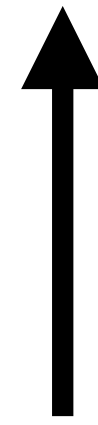
Overview

- Traditional speech recognition
- Self-supervised learning for speech processing
 - wav2vec 2.0
 - Cross-lingual training
 - Completely unsupervised speech recognition

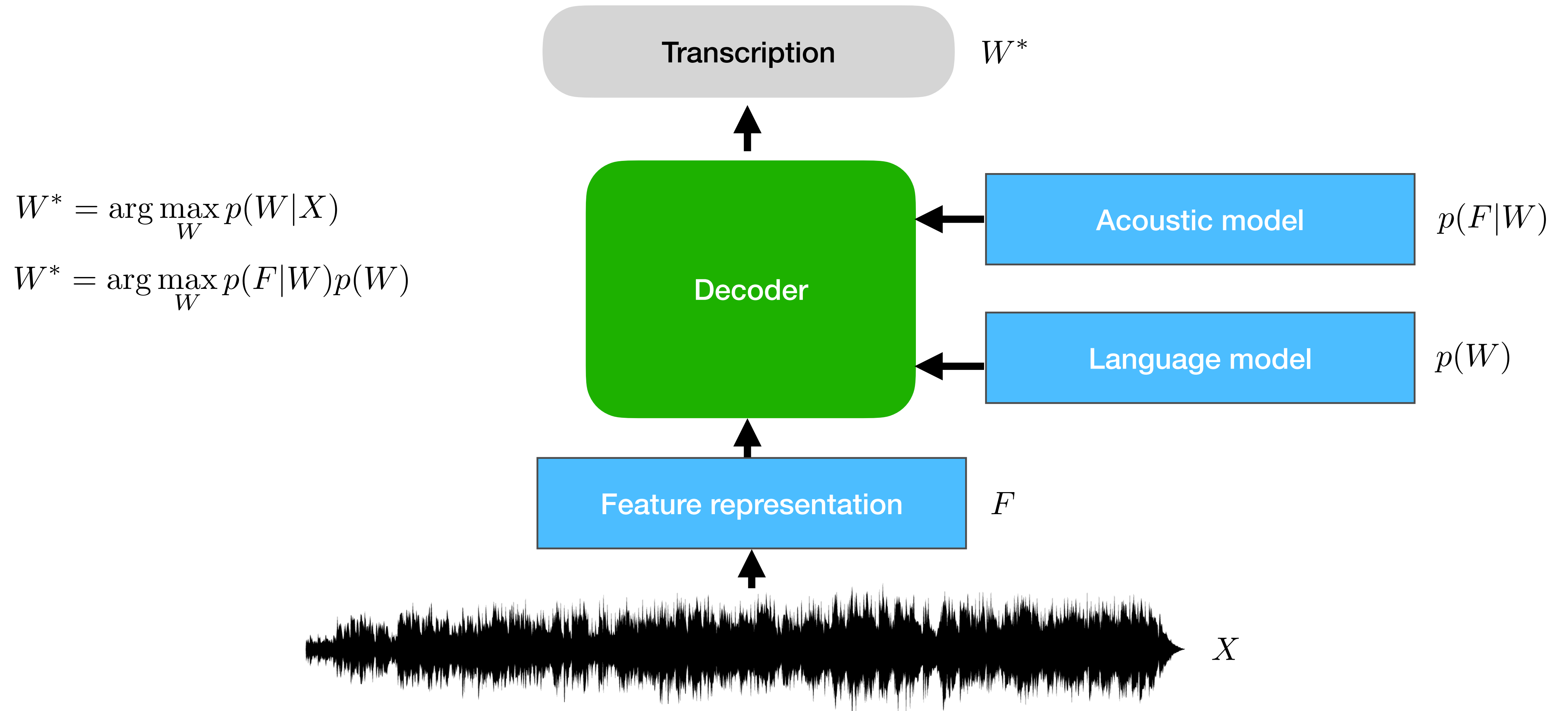
Traditional speech recognition

Speech recognition

I like black tea with milk



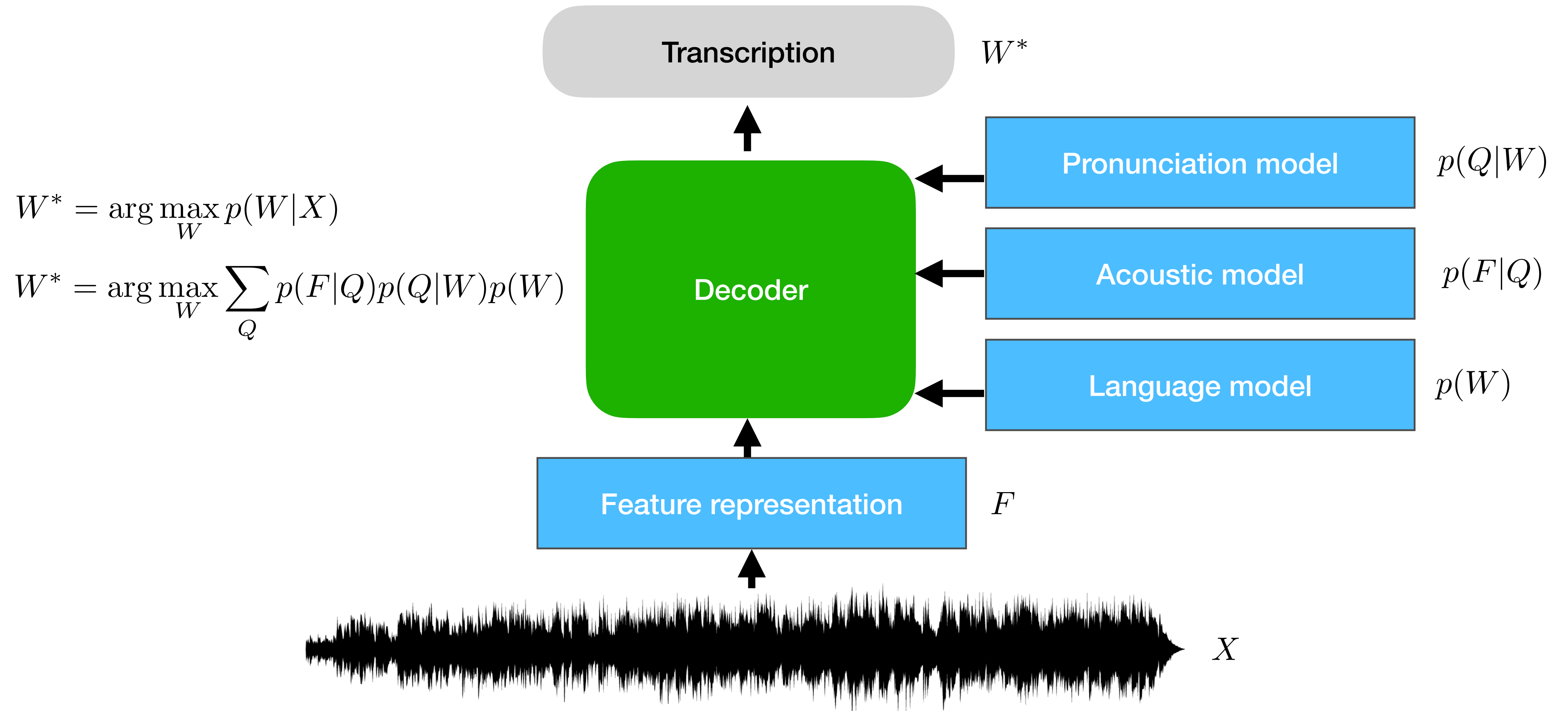
Traditional automatic speech recognition (ASR)



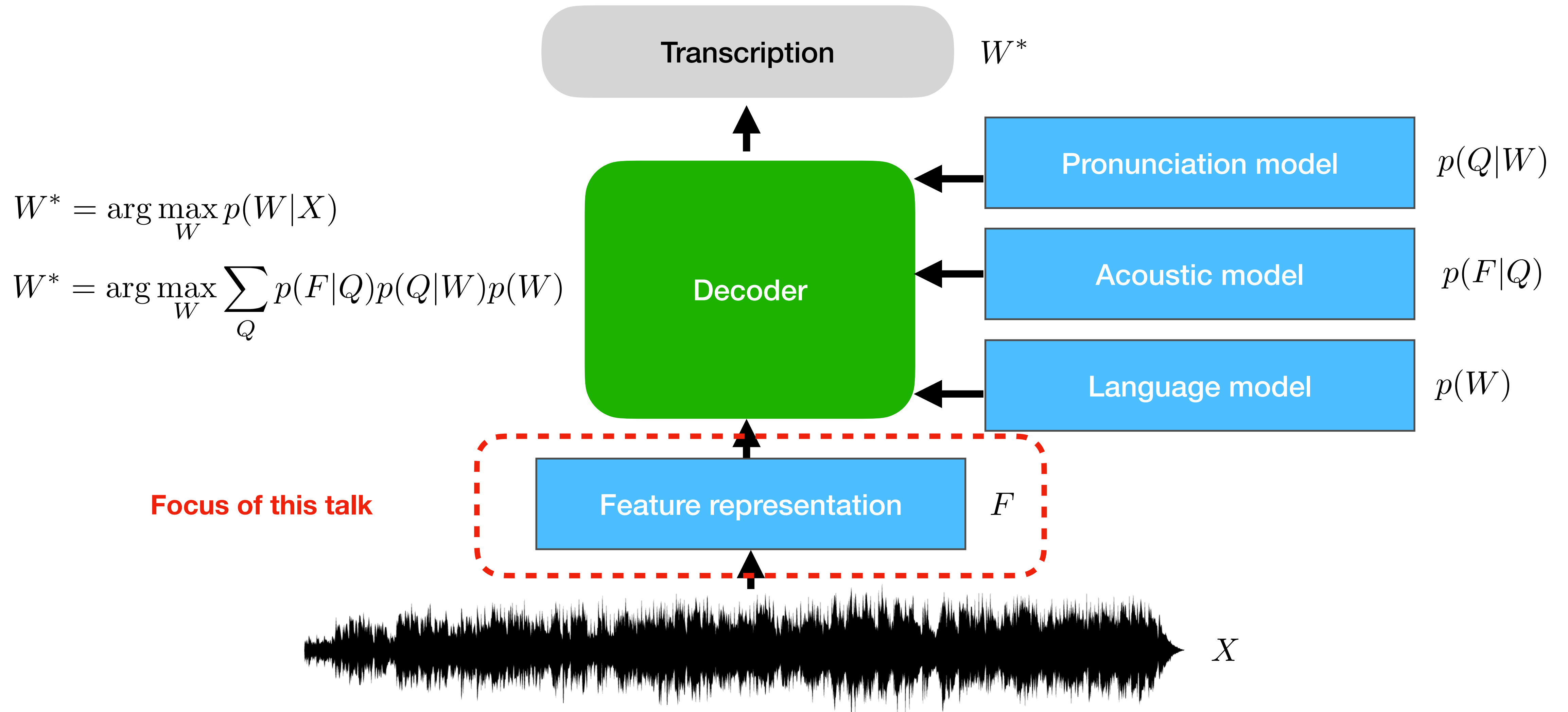
Traditional automatic speech recognition (ASR)

- Represent words as sequences of phonemes
- hello = h eh l ow
- Distinct units of sound to distinguish words

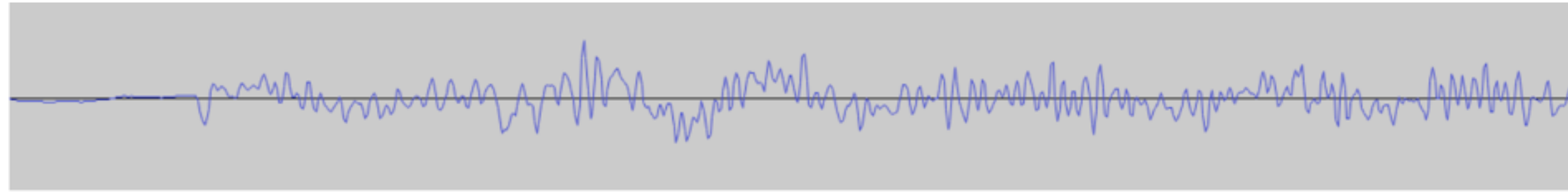
Traditional automatic speech recognition (ASR)



Traditional automatic speech recognition (ASR)



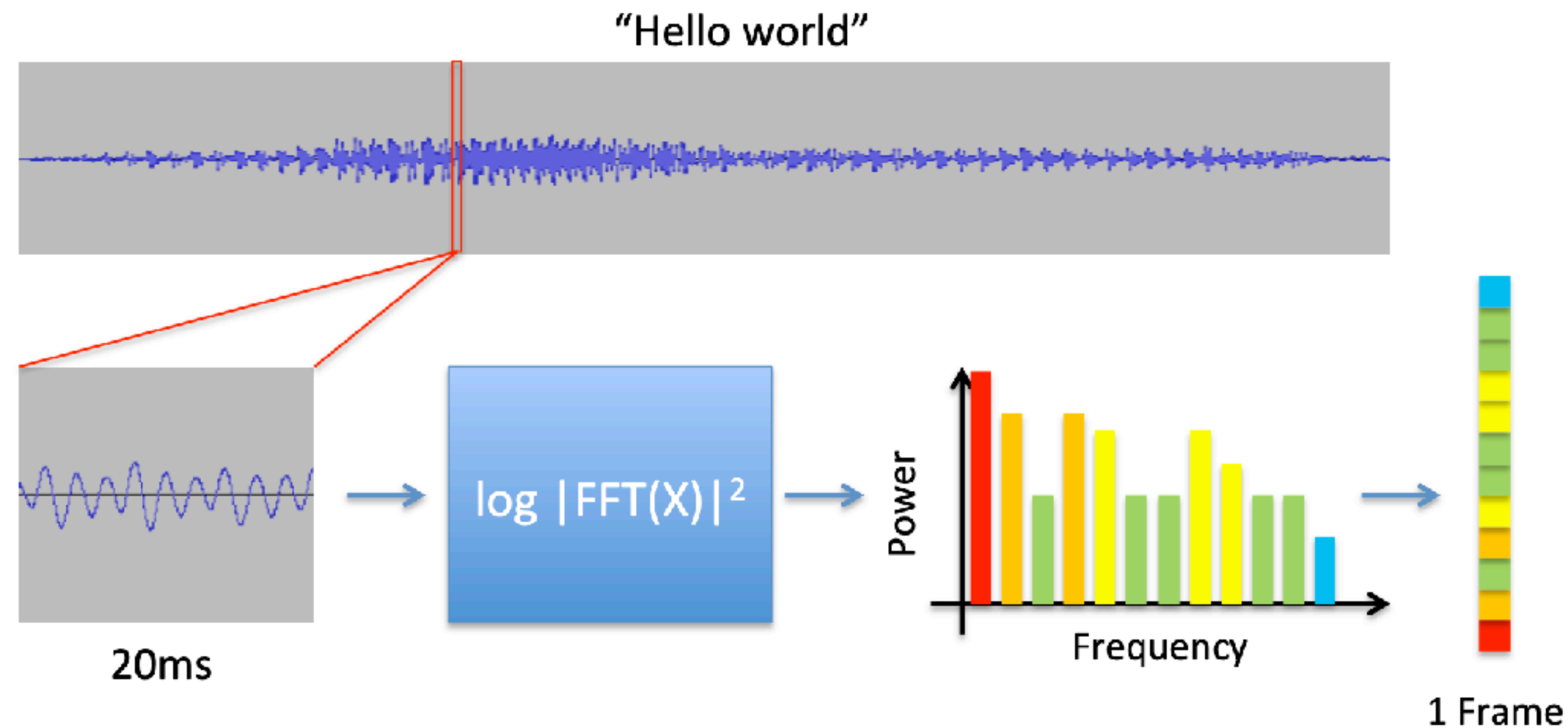
Feature representation



- Typical sample rates for speech: 8KHz, 16KHz.
- Traditionally: build spectrogram

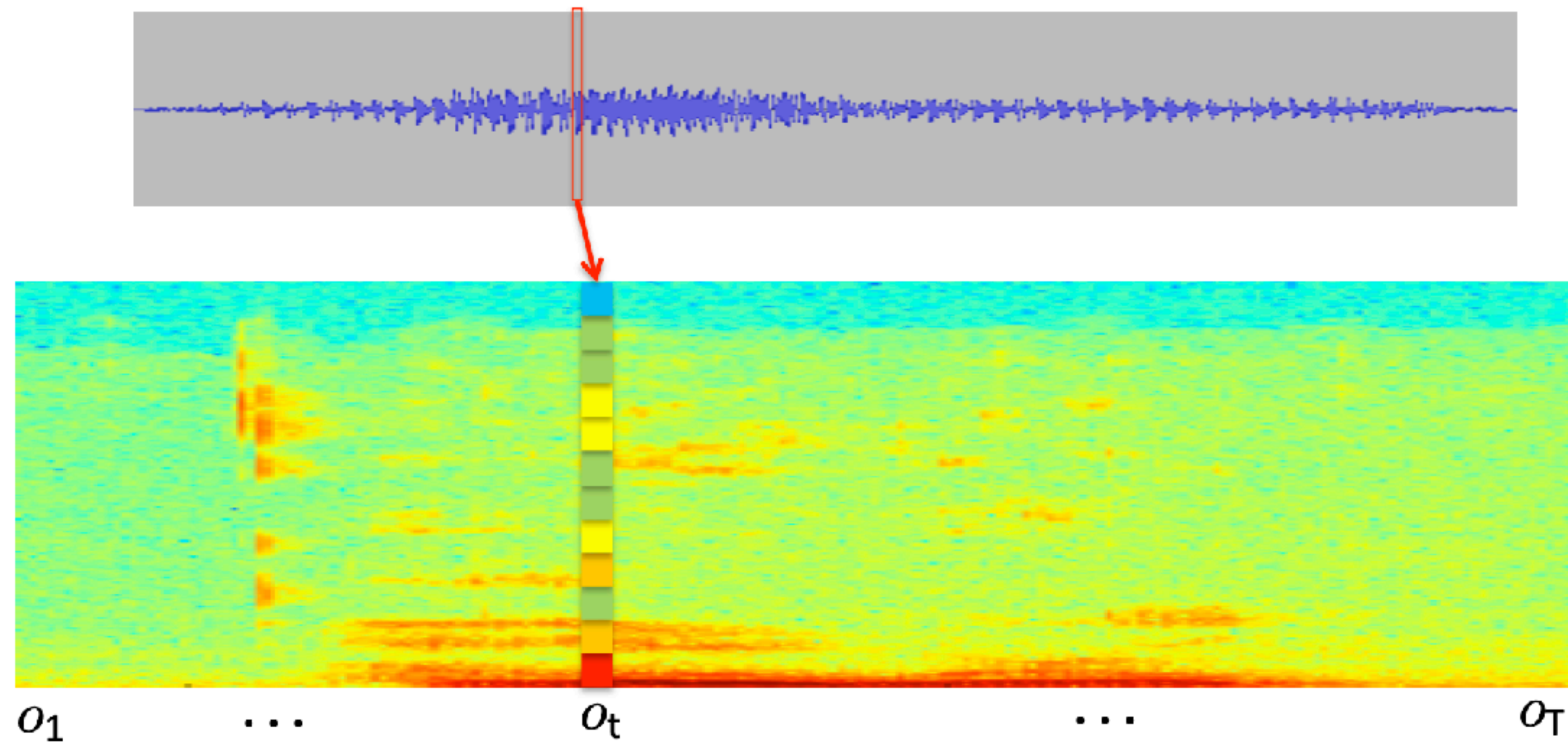
Spectrogram

- Small window, e.g., 20ms of waveform
- Compute FFT and take magnitude
- Describes frequency content in local window



Spectrogram

- Concatenate frames from adjacent windows to form a spectrogram





Self-supervised speech representation learning

Training speech recognition models

I like black tea with milk



- Train on 1,000s of hours of transcribed data for good systems.
- Many languages, dialects, domains etc.



Supervised machine learning



potential train/test mismatch



Need to annotate lots of data!

Supervised machine learning

( , cat)

Not how humans learn!

potential train/test mismatch

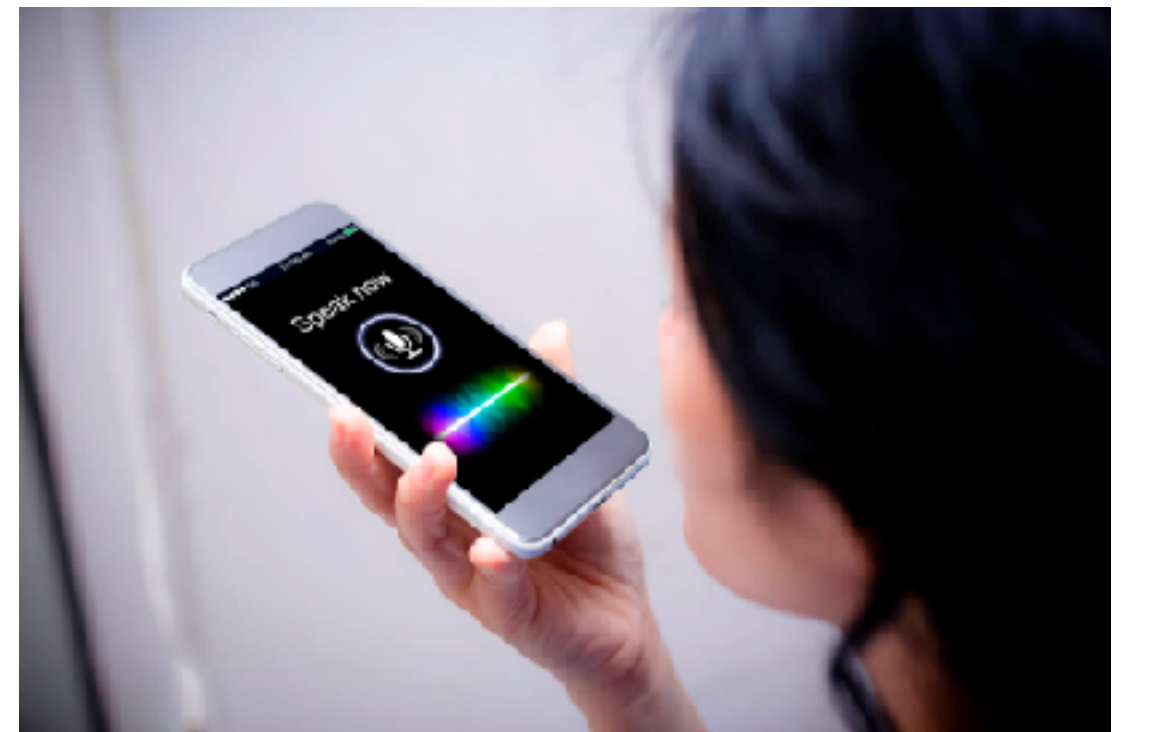
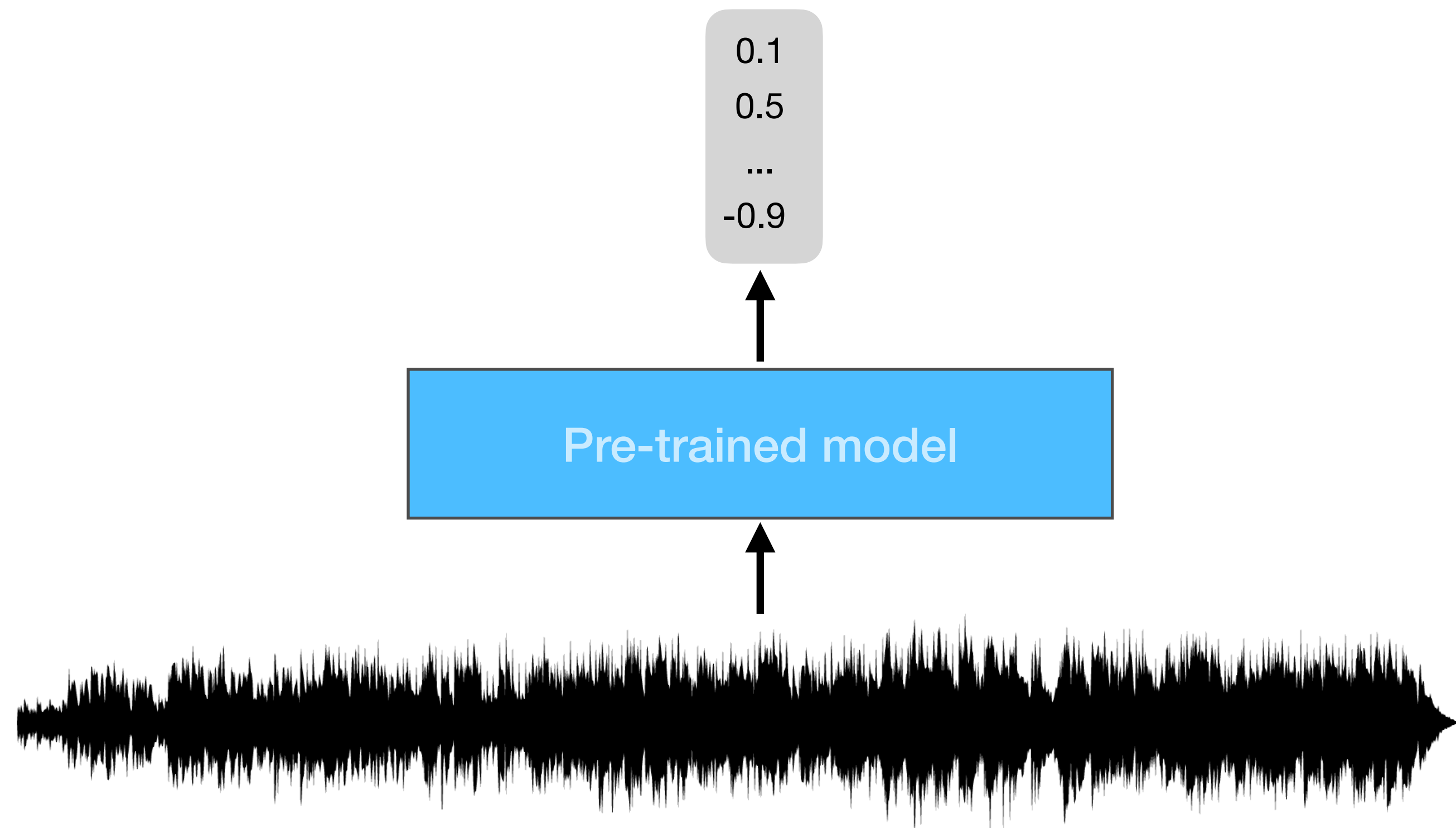


Need to annotate lots of data!

Supervised machine learning



Learning good representations of audio data
from unlabeled audio

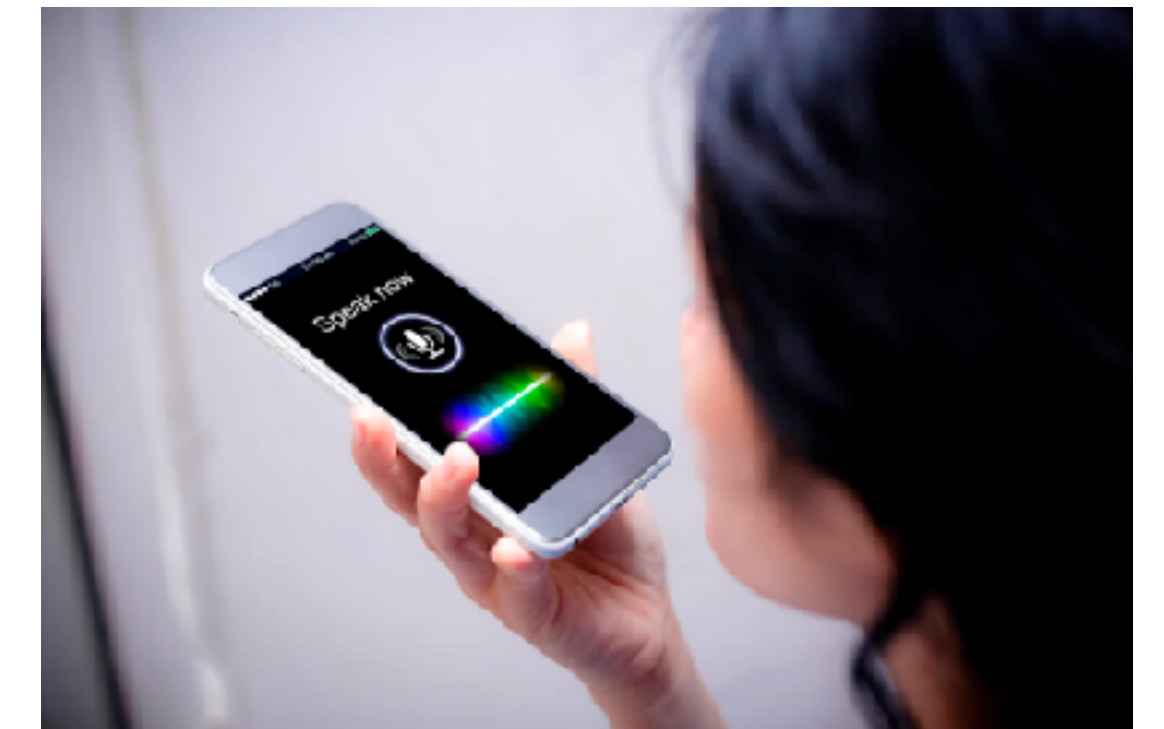
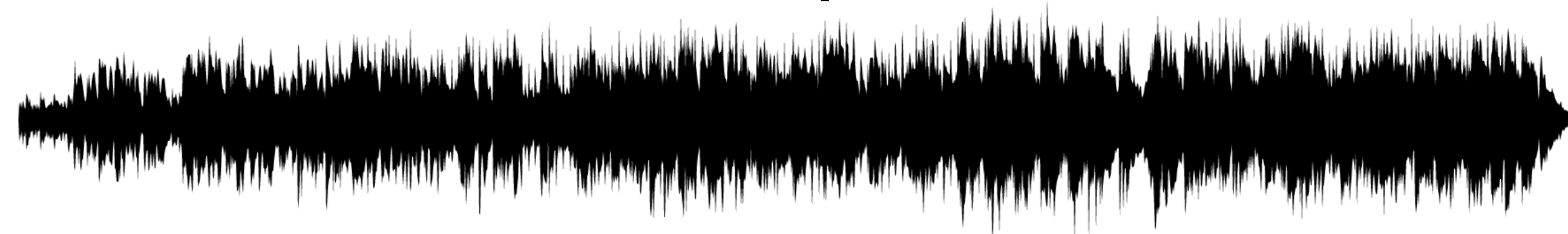


I like tea

Speech recognition

0.1
0.5
...
-0.9

Pre-trained model



Speech translation

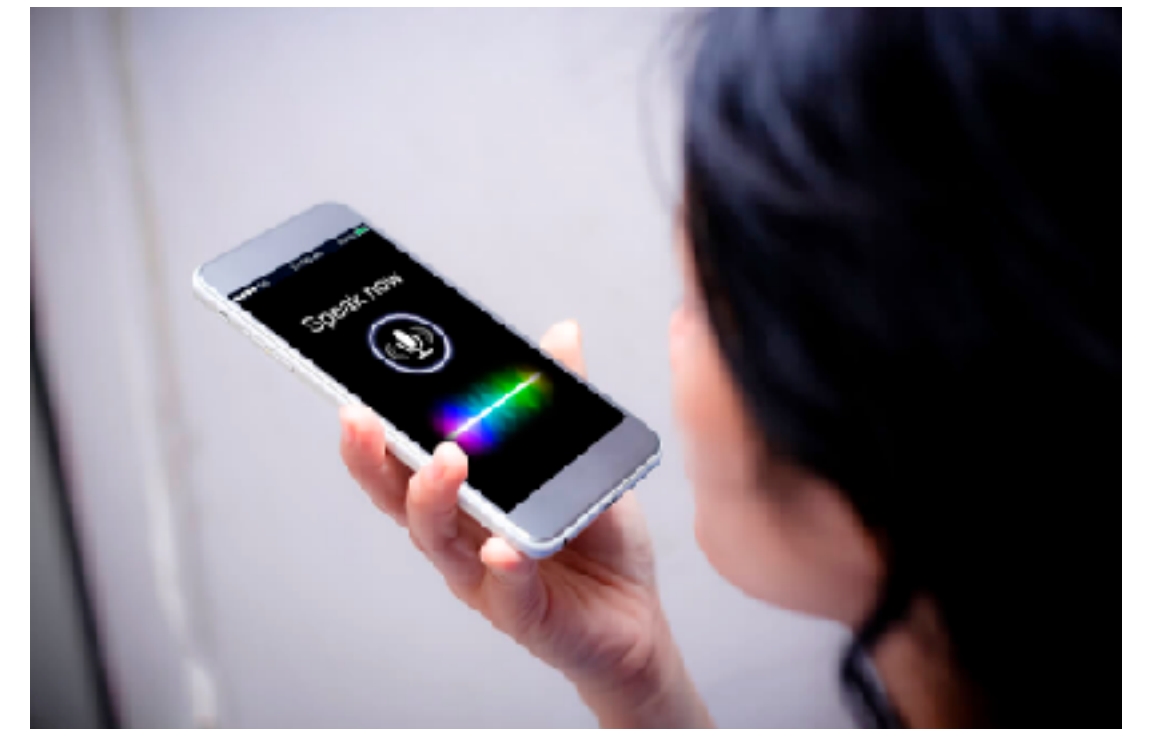


Pre-trained model

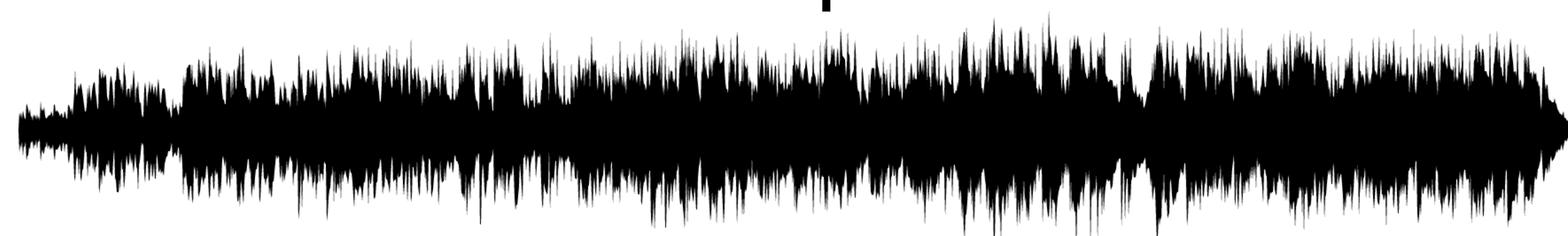
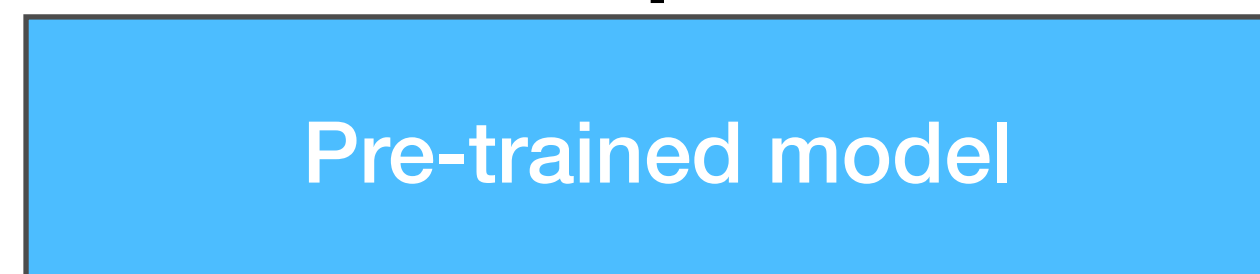
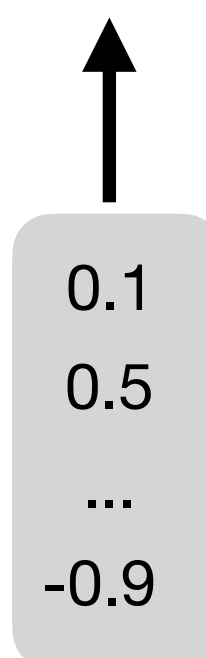
0.1
0.5
...
-0.9



Ich mag Tee



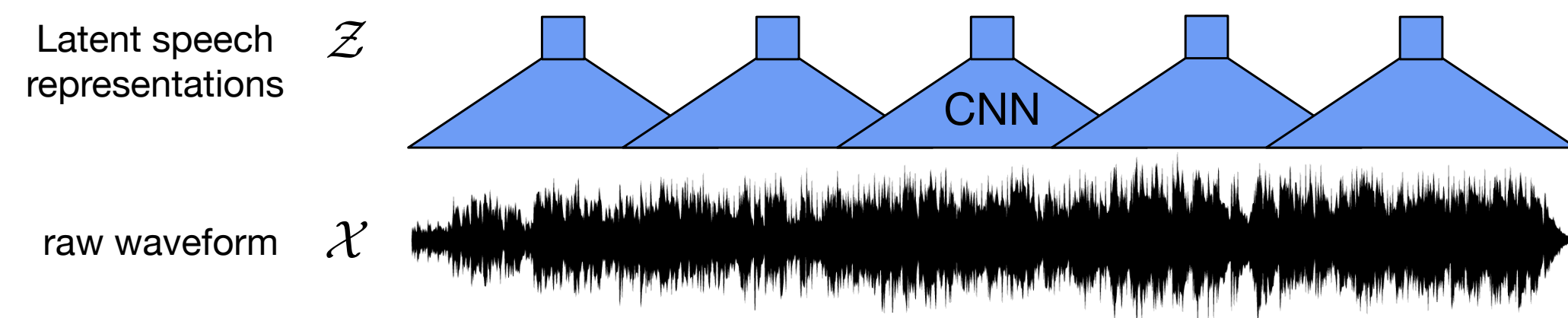
"music"



Audio event detection

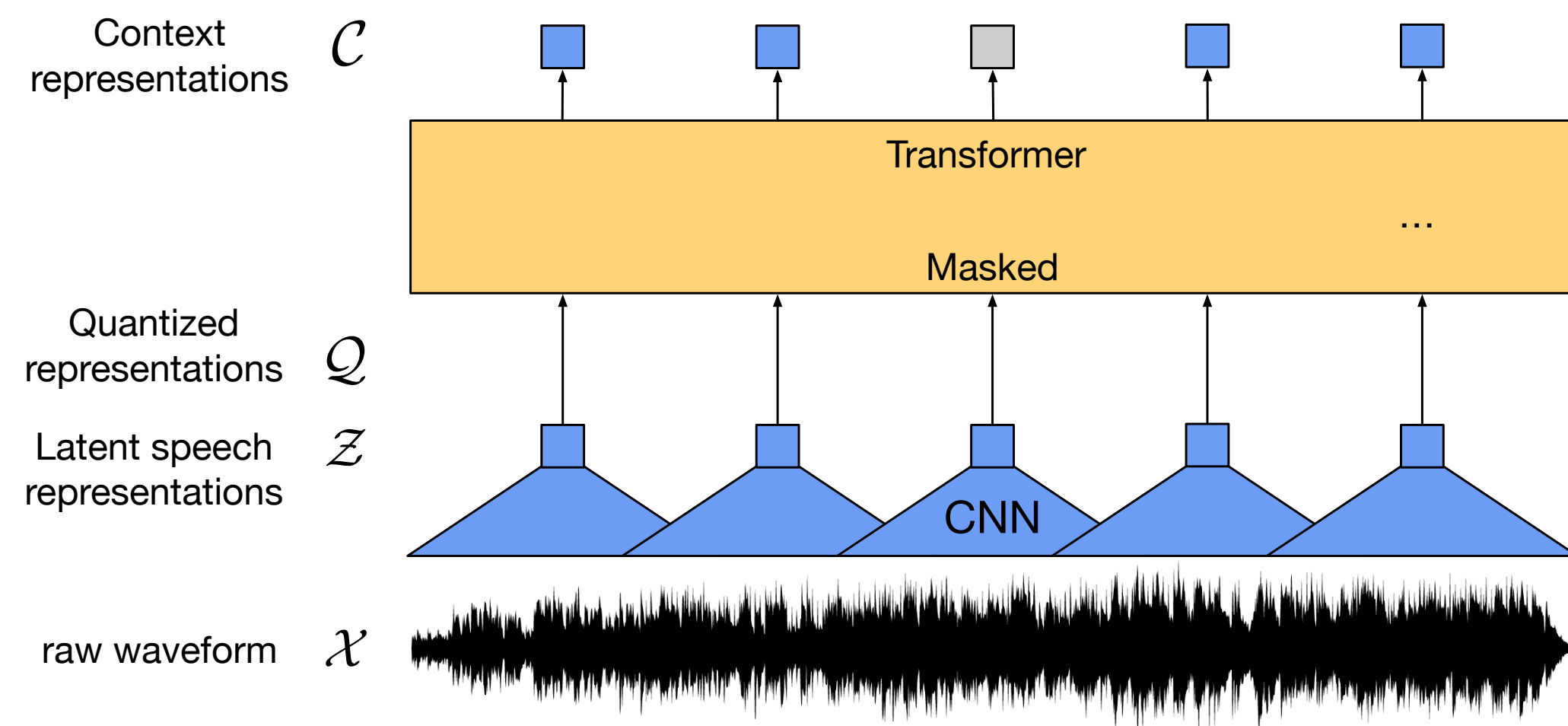


wav2vec 2.0



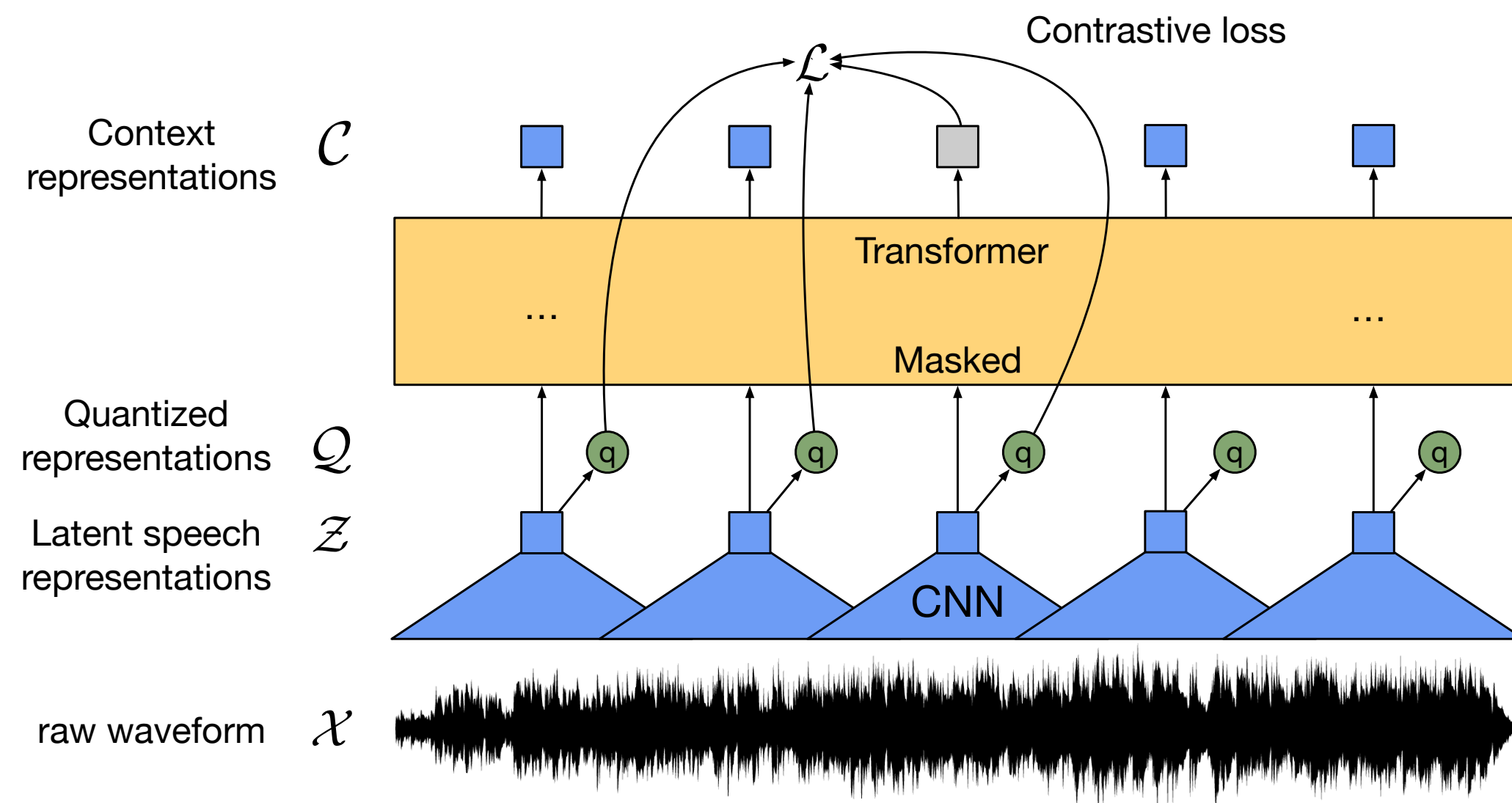
- Masked prediction with transformer, bi-directional contextualized representations (similar to BERT).
- But predict what? Learn an inventory of speech units with vector quantization via Gumbel softmax.
- Learning task: Joint VQ & context representation learning.
- Contrast true quantized latent with distractor latents.

wav2vec 2.0



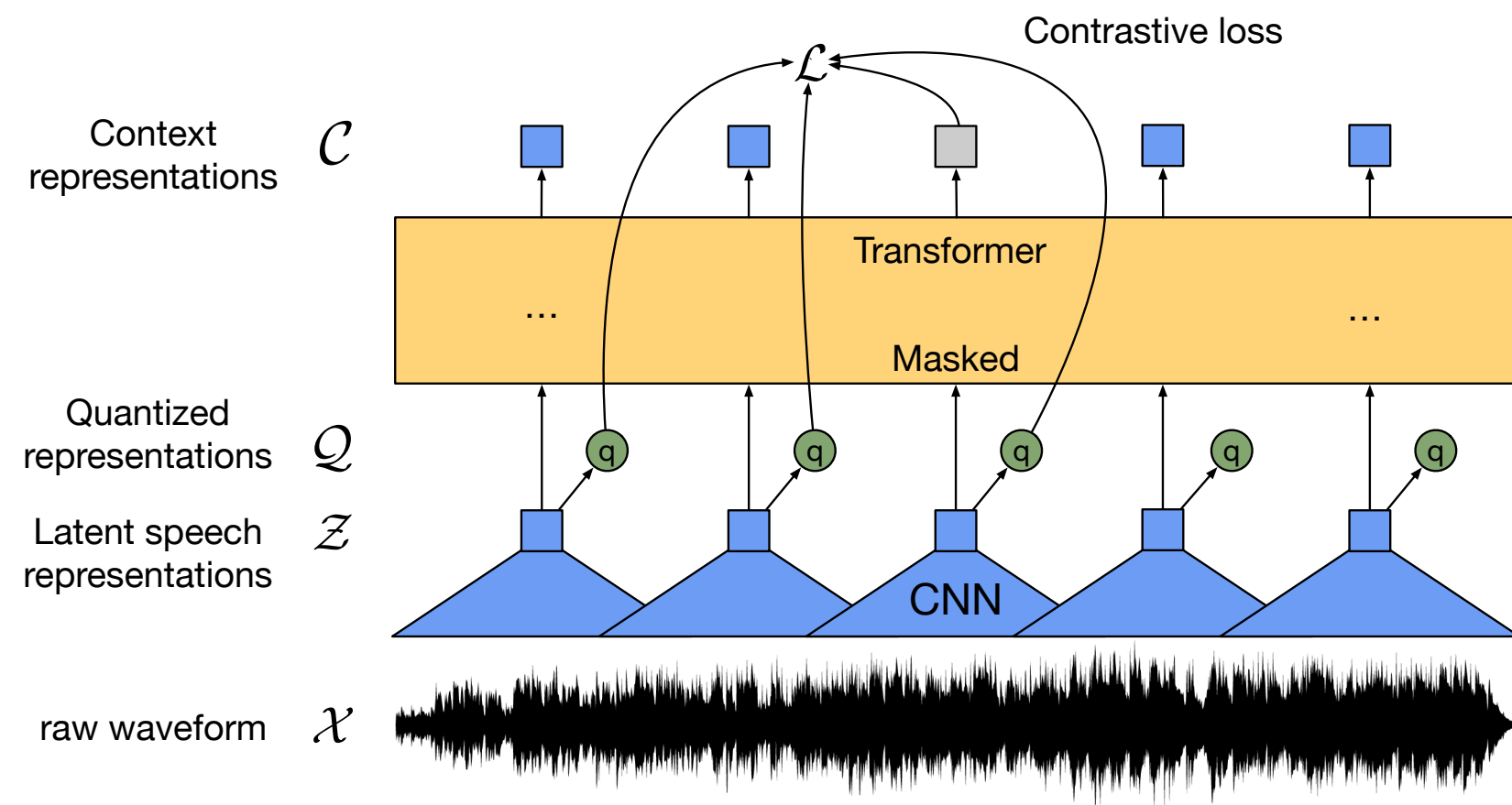
- Masked prediction with transformer, bi-directional contextualized representations (similar to BERT).
- But predict what? Learn an inventory of speech units with vector quantization via Gumbel softmax.
- Learning task: Joint VQ & context representation learning.
- Contrast true quantized latent with distractor latents.

wav2vec 2.0



- Masked prediction with transformer, bi-directional contextualized representations (similar to BERT).
- But predict what? Learn an inventory of speech units with vector quantization via Gumbel softmax.
- Learning task: Joint VQ & context representation learning.
- Contrast true quantized latent with distractor latents.

Objective



Cosine similarity

Context representation

Discrete latent speech representation

$$\mathcal{L}_m = -\log \frac{\exp(\text{sim}(\mathbf{c}_t, \mathbf{q}_t)/\kappa)}{\sum_{\tilde{\mathbf{q}} \sim \mathbf{Q}_t} \exp(\text{sim}(\mathbf{c}_t, \tilde{\mathbf{q}})/\kappa)}$$

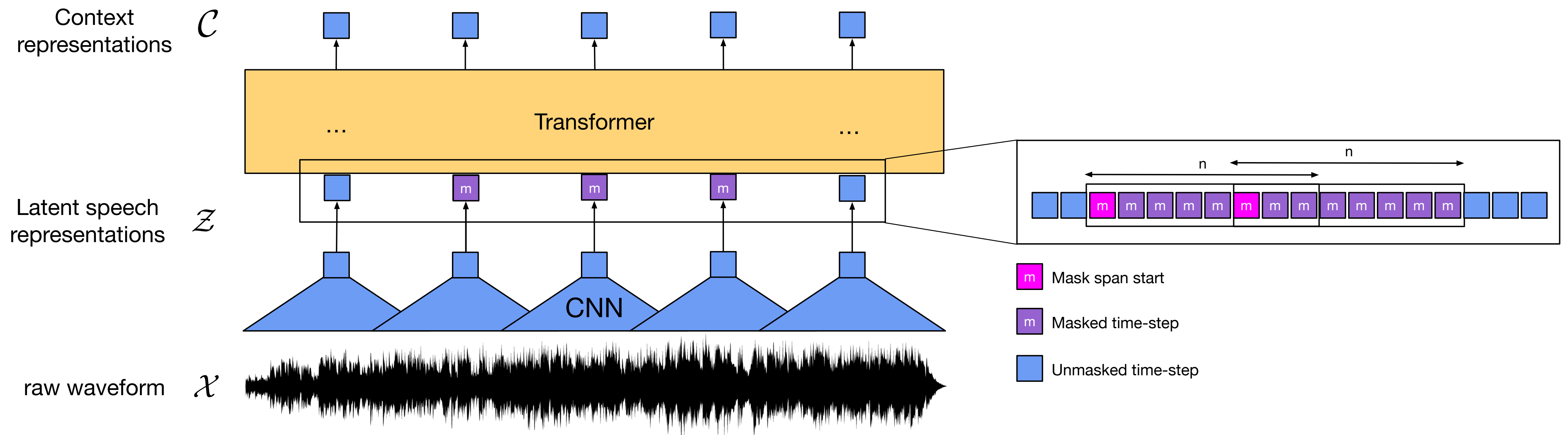
Negative samples

Temperature

Codebook diversity penalty to encourage more codes to be used

Masking

- Sample starting points for masks without replacement, then expand to 10 time-steps (1 time-step is 25ms but 10ms stride)
- Spans can overlap
- For a 15s sample, ~49% of the time-steps masked with an average span length of ~300ms

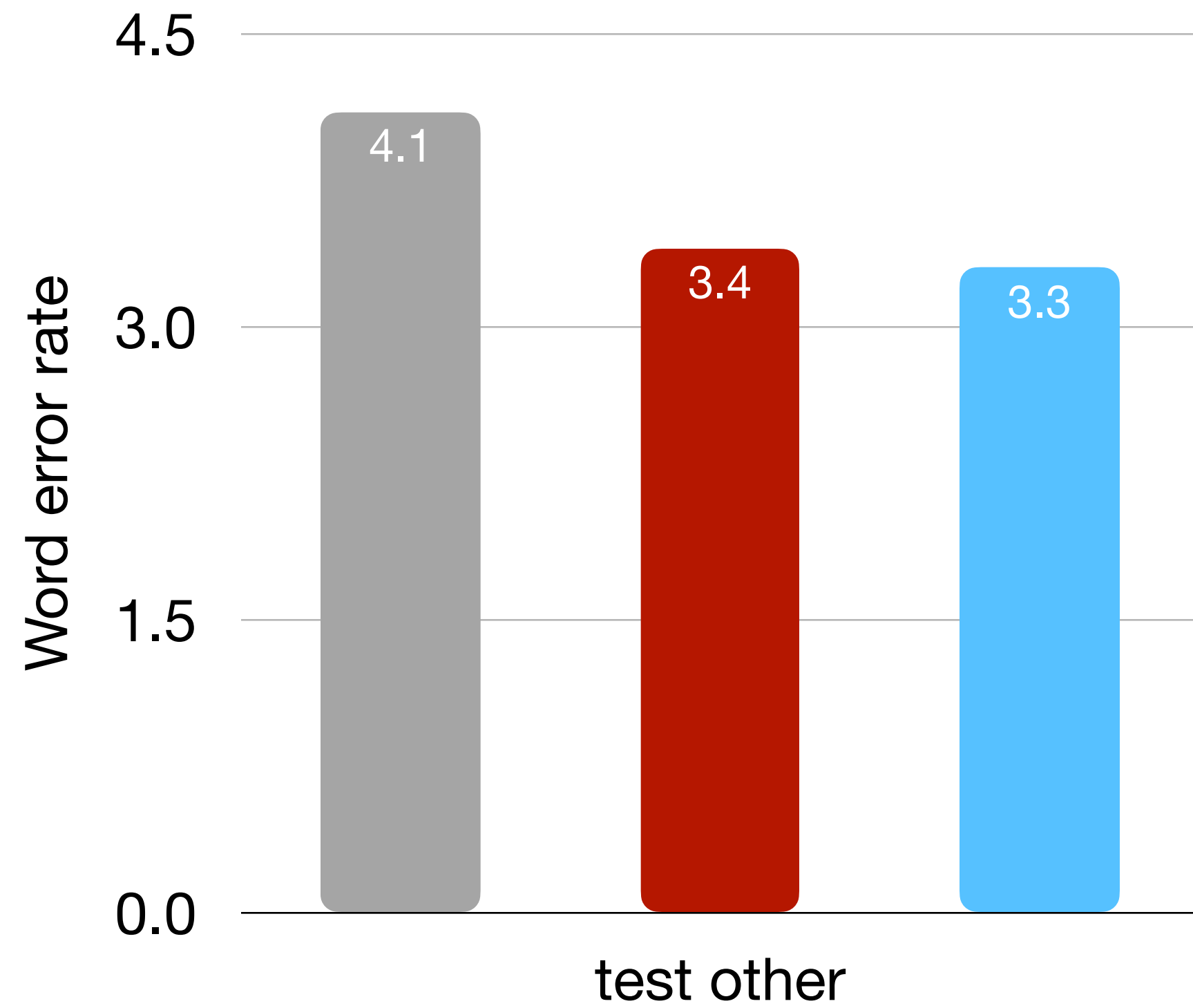


Fine-tuning

- Add a single linear projection on top into target vocab and train with CTC loss with a low learning rate (CNN encoder is not trained).
- Use modified SpecAugment in latent space to prevent early overfitting
- Uses wav2letter decoder with the official 4gram LM and Transformer LM

Results

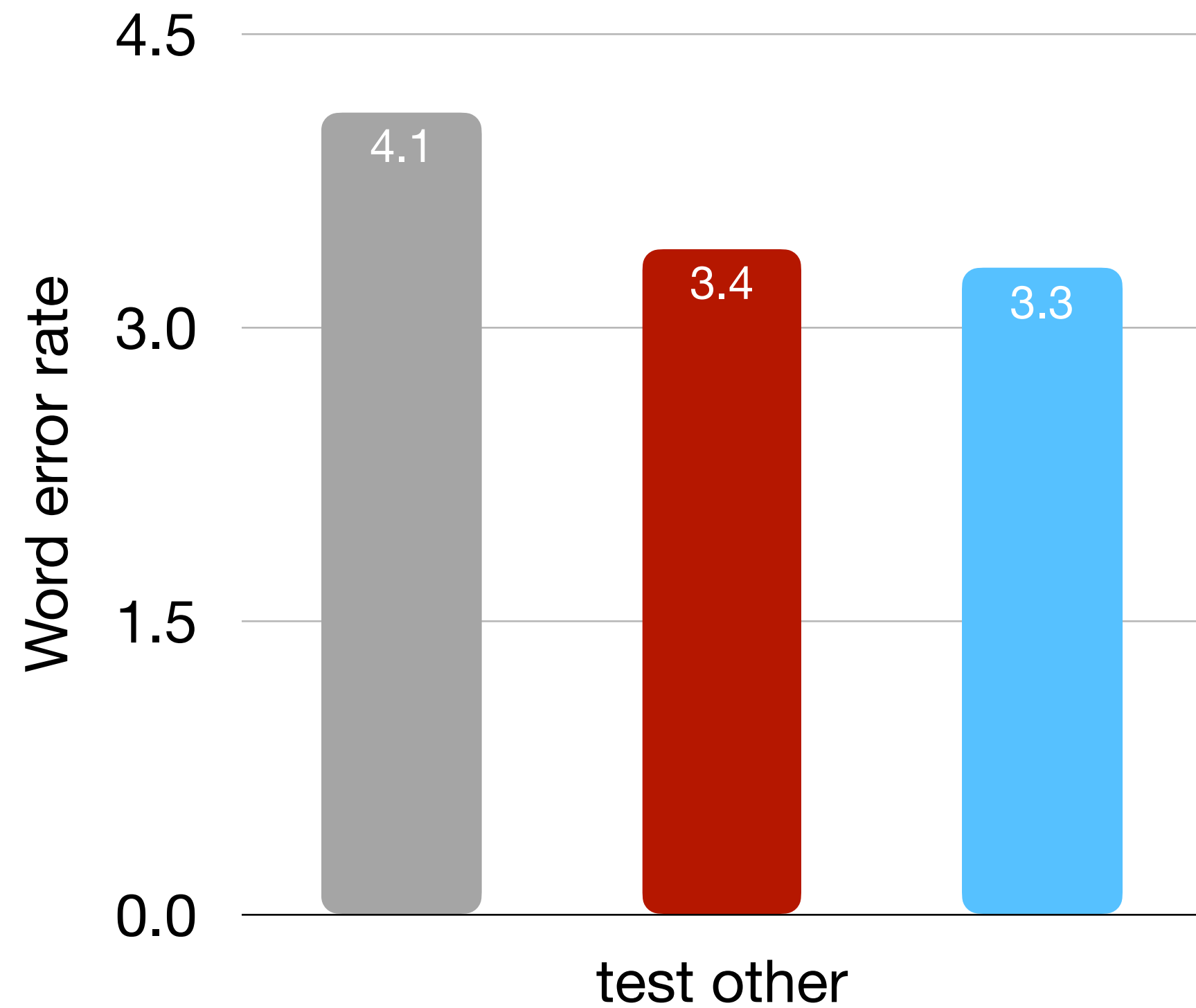
High resource
(Librispeech 960h labeled)



- ContextNet (supervised)
- Noisy Student (60k-h unlabeled)
- wav2vec (60k-h unlabeled)

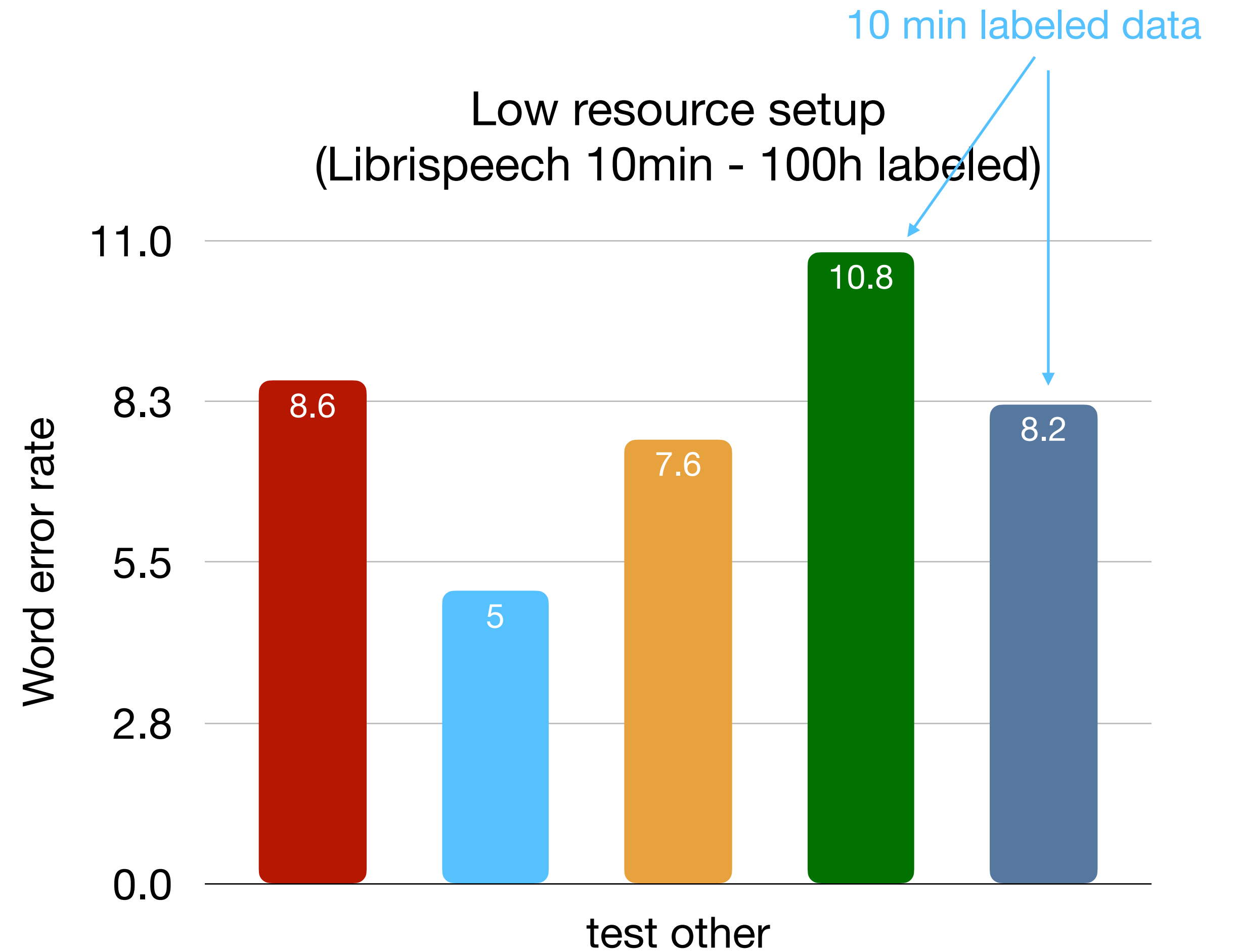
Results

High resource
(Librispeech 960h labeled)



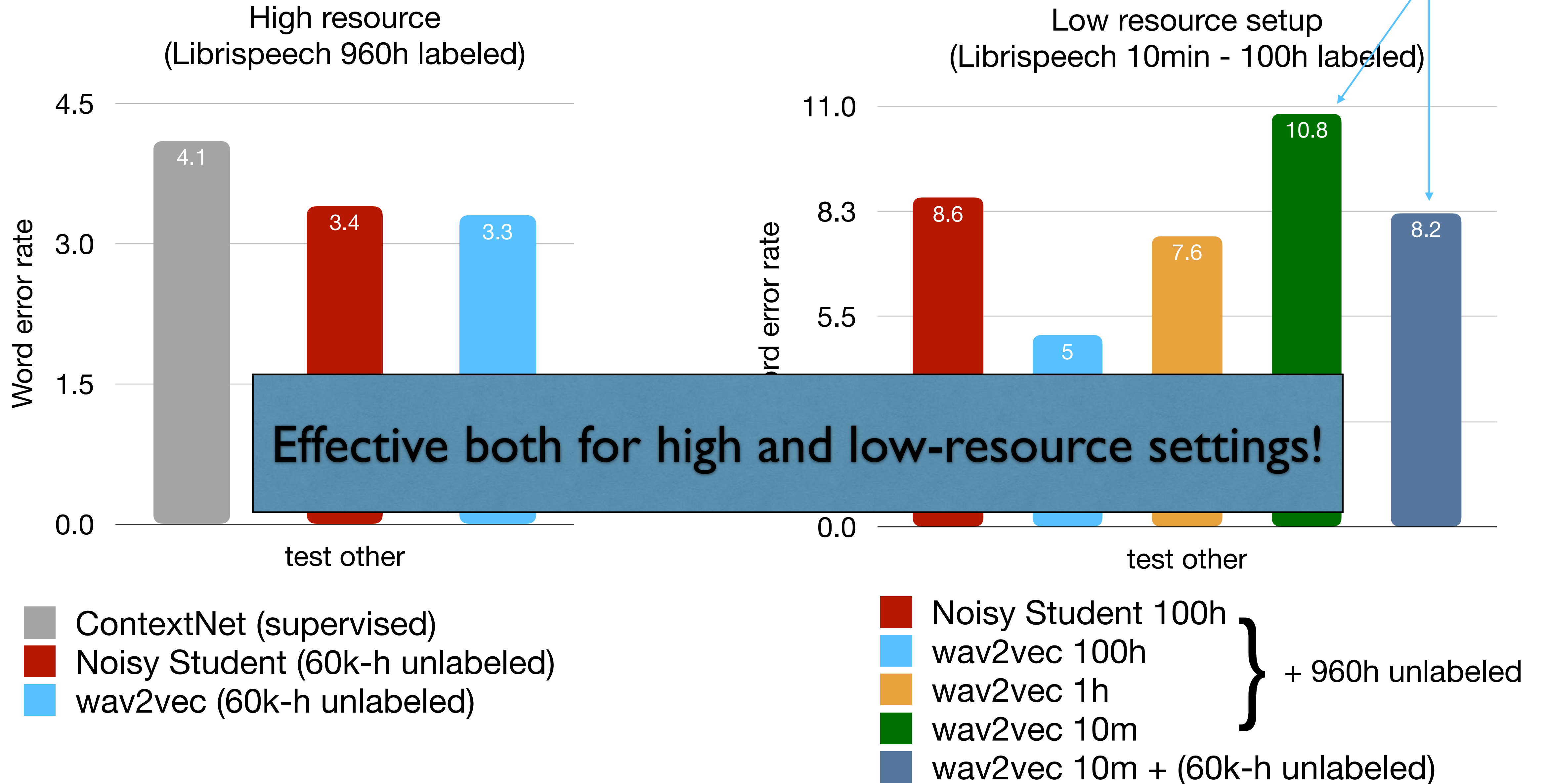
- ContextNet (supervised)
- Noisy Student (60k-h unlabeled)
- wav2vec (60k-h unlabeled)

Low resource setup
(Librispeech 10min - 100h labeled)

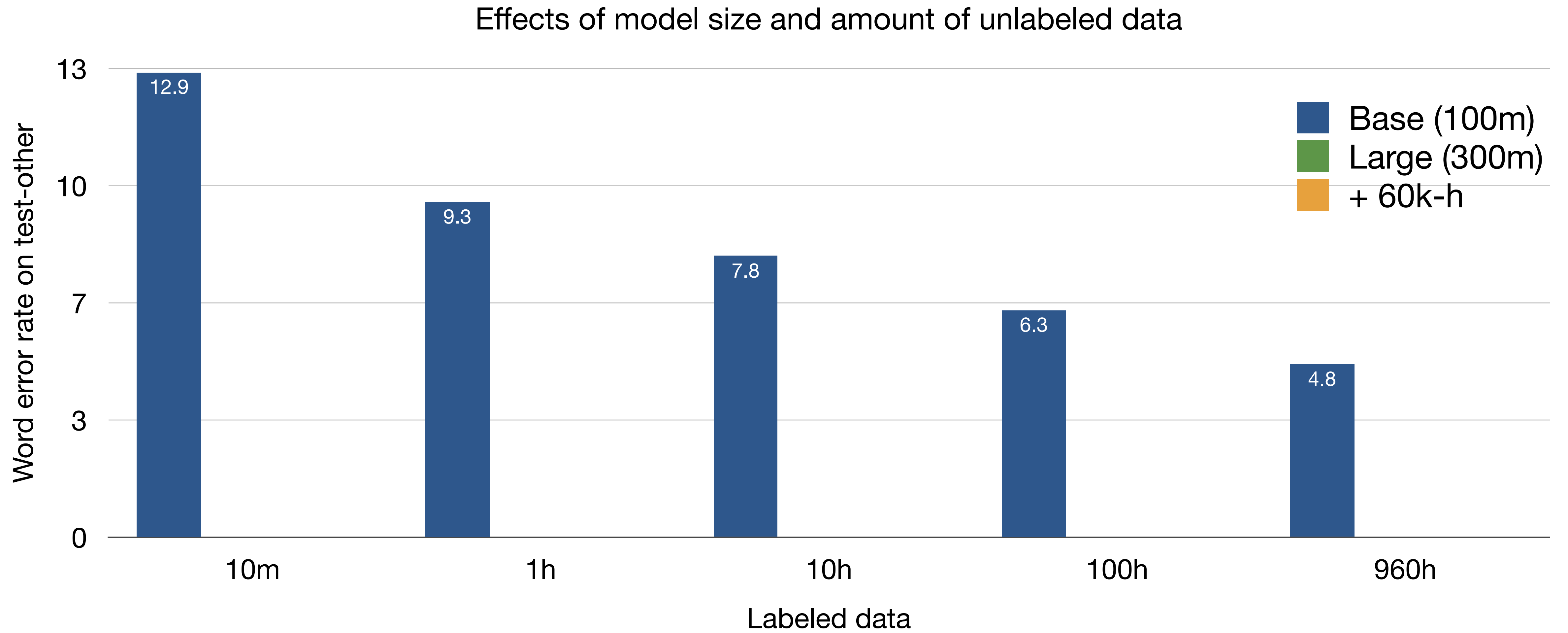


- Noisy Student 100h
 - wav2vec 100h
 - wav2vec 1h
 - wav2vec 10m
 - wav2vec 10m + (60k-h unlabeled)
- } + 960h unlabeled

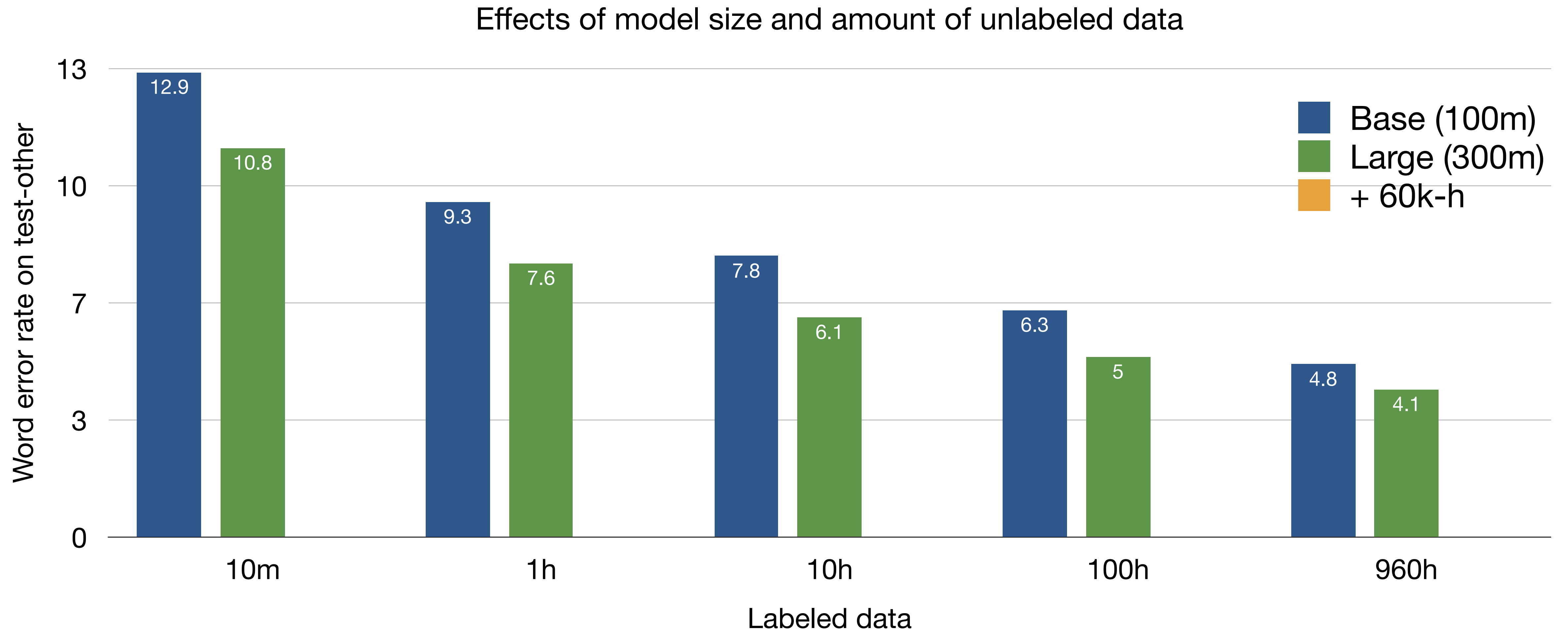
Results



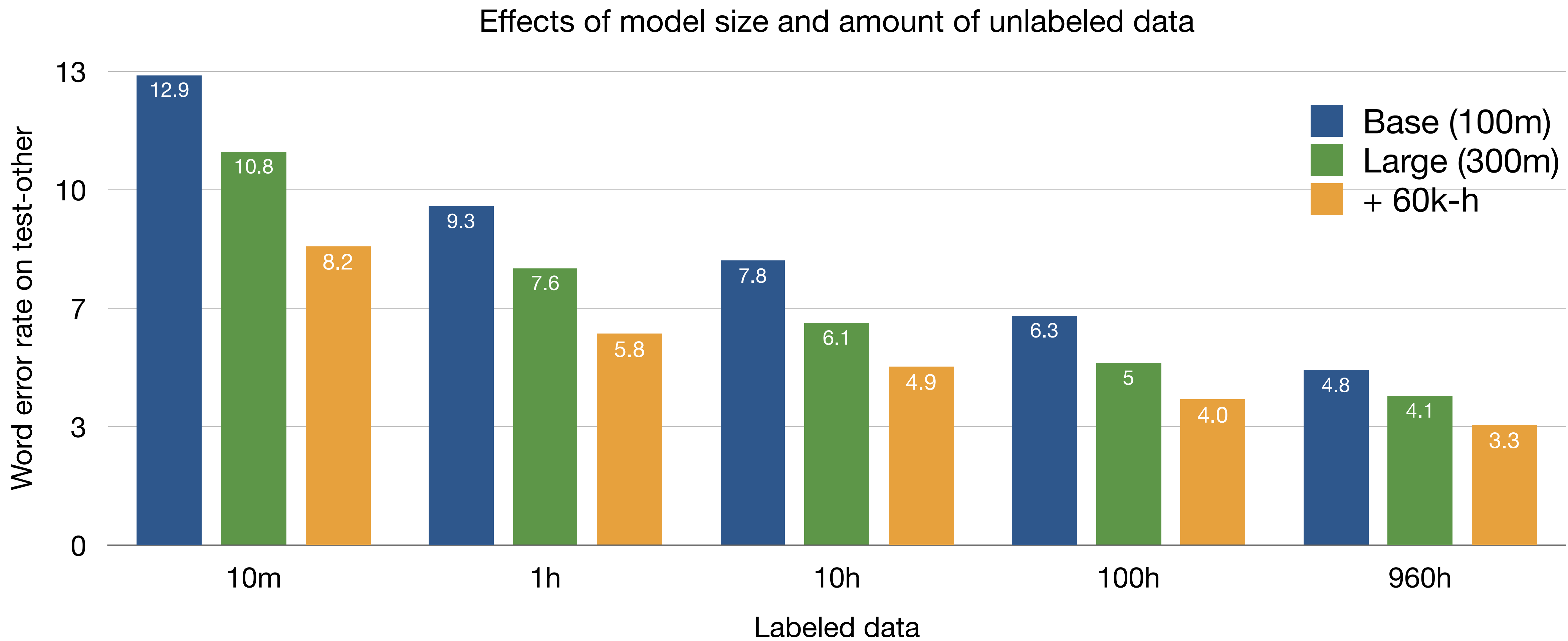
Results



Results



Results



Examples (10 min labeled data)

HYP (no LM): she SESED and LUCHMAN GAIVE A SENT won by her GENTAL argument

HYP (w/ LM): she ceased and LUCAN gave assent won by her gentle argument

REF: she ceased and lakshman gave assent won by her gentle argument

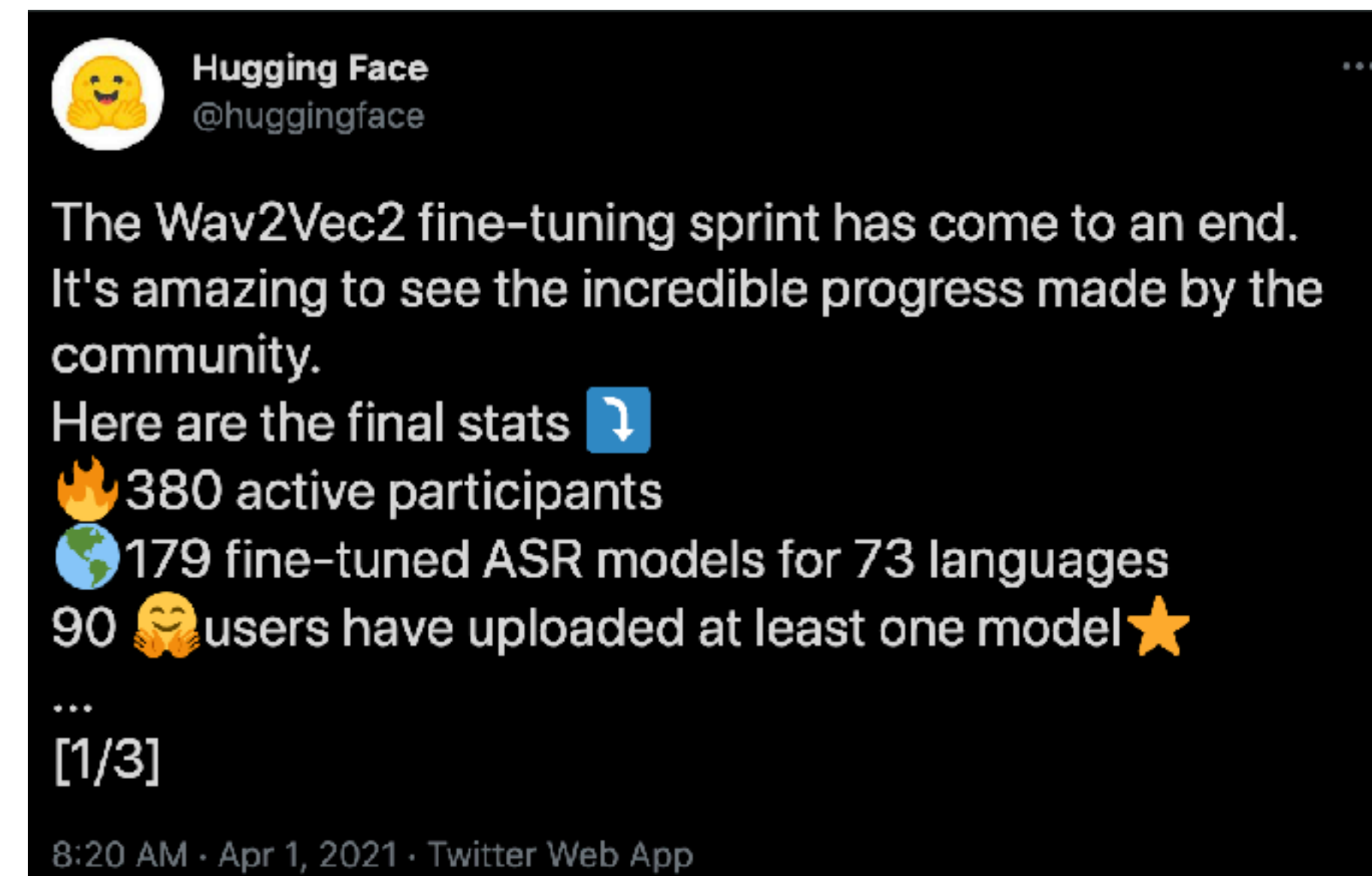
HYP (no LM): but NOT WITH STANDING this boris EMBRAED him in a QUIAT FRIENDLY way and CISED him THREE times

HYP (w/ LM): but NOT WITHSTANDING this boris embraced him in a quiet friendly way and kissed him three times

REF: but notwithstanding this boris embraced him in a quiet friendly way and kissed him three times

wav2vec on Hugging Face

- Hugging Face is a popular NLP model zoo
- Hugging Face community fine-tuned our models to do speech recognition in 73 languages.



Pre-training and self-training

- Self-training very successful in speech recognition: generate pseudo-labels



Supervised model

Pre-training and self-training

- Self-training very successful in speech recognition: generate pseudo-labels

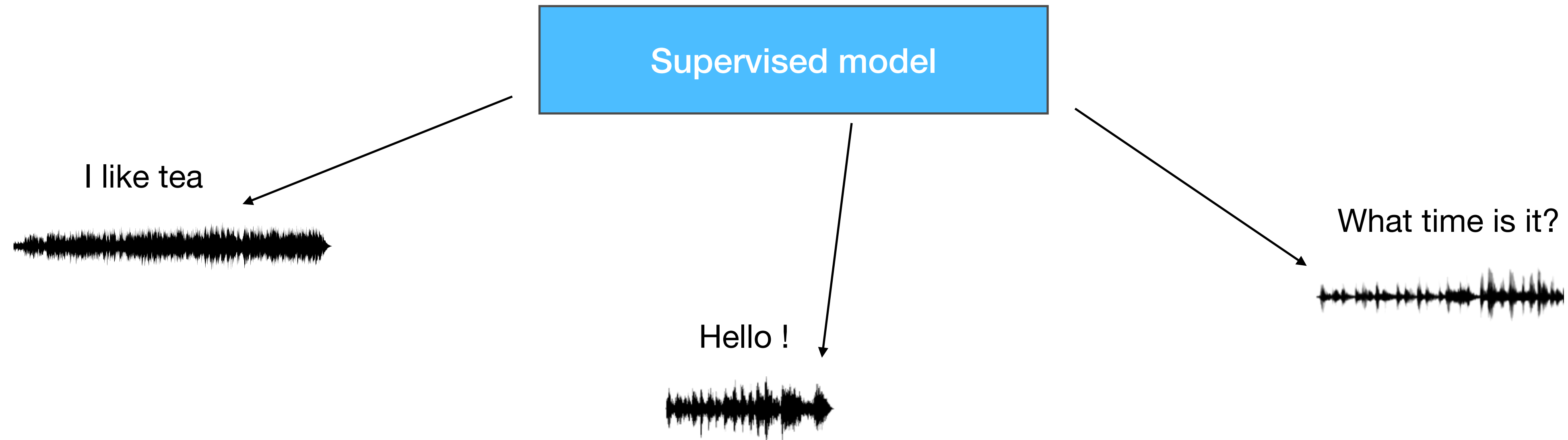


Supervised model

The diagram illustrates a supervised model's output. A blue rectangular box labeled "Supervised model" is positioned at the top center. Below it, three black audio waveforms are arranged horizontally. The first waveform on the left is a dense, high-frequency signal. The second waveform in the middle is a shorter, lower-frequency signal. The third waveform on the right is a signal with a distinct, periodic pattern, resembling a vowel or a specific phoneme.

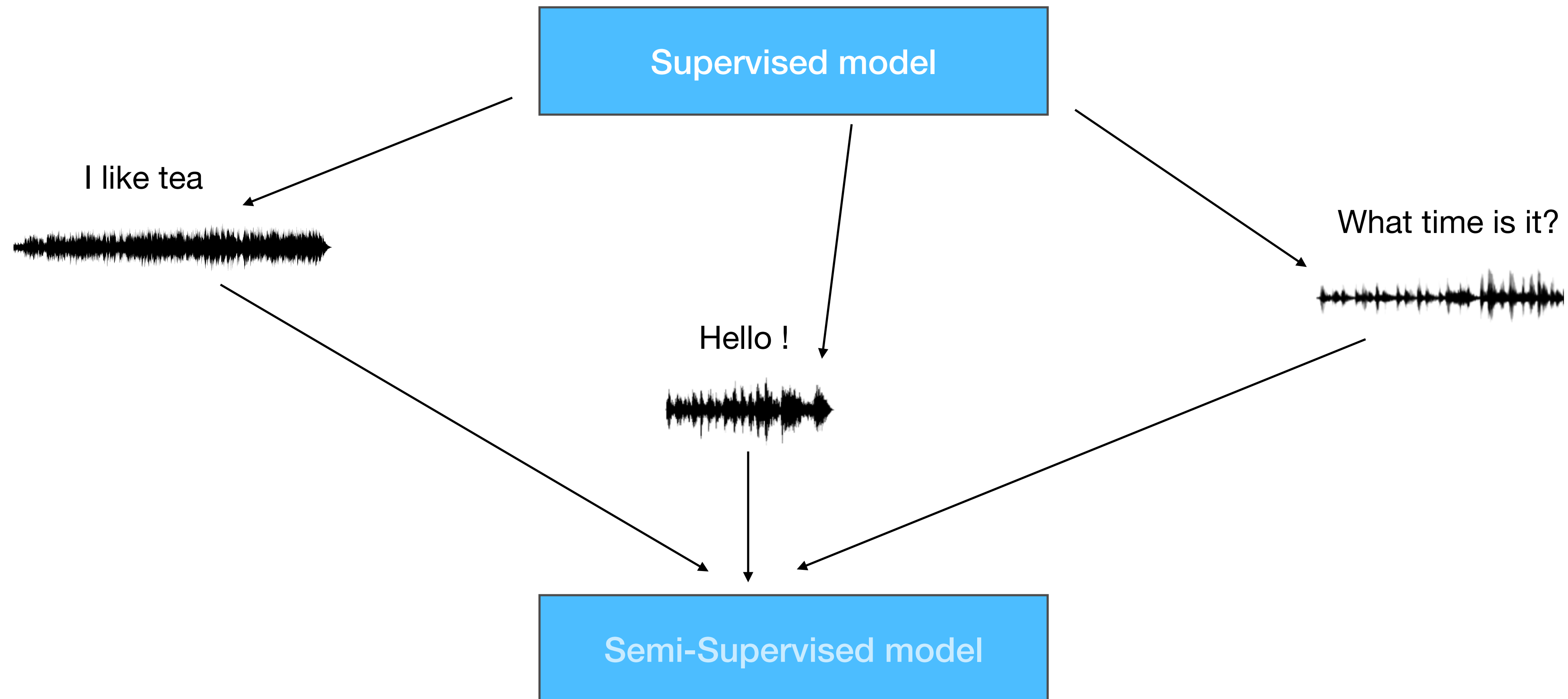
Pre-training and self-training

- Self-training very successful in speech recognition: generate pseudo-labels

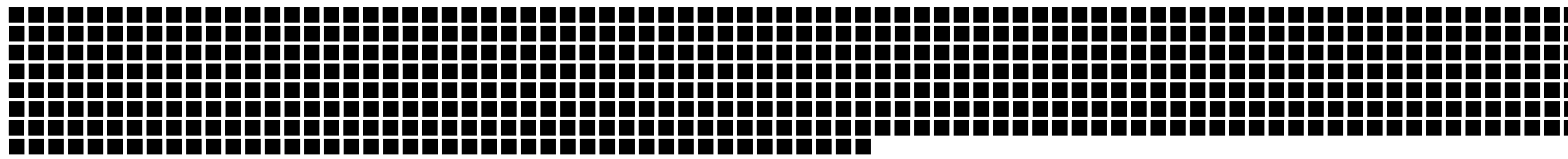


Pre-training and self-training

- Self-training very successful in speech recognition: generate pseudo-labels

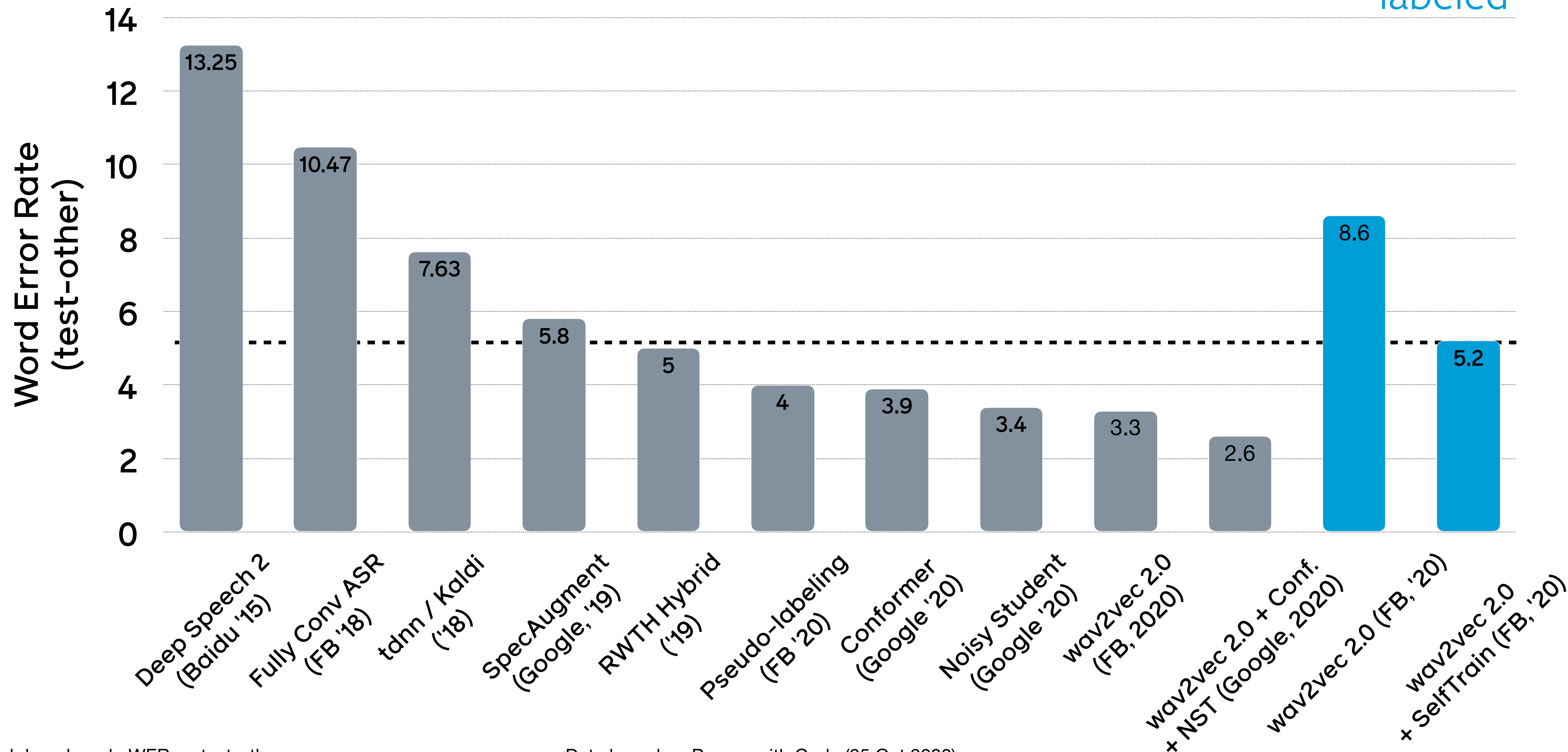


Amount of
labeled
data used

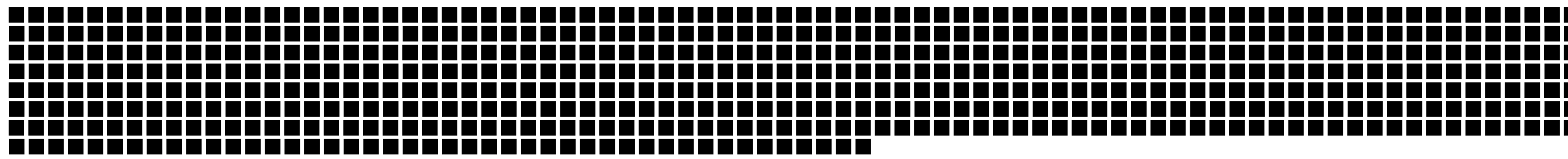


960h labeled

↑
10min
labeled

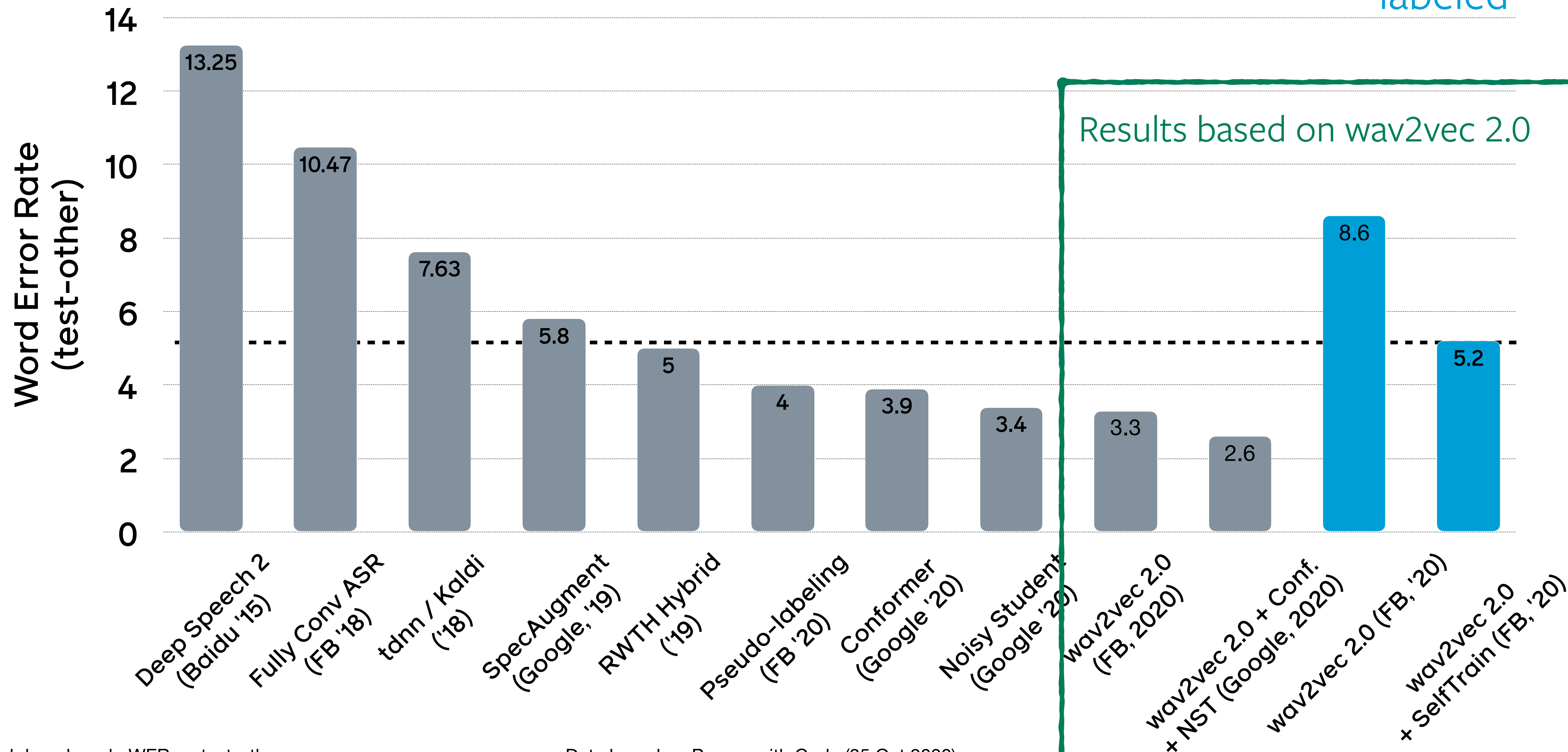


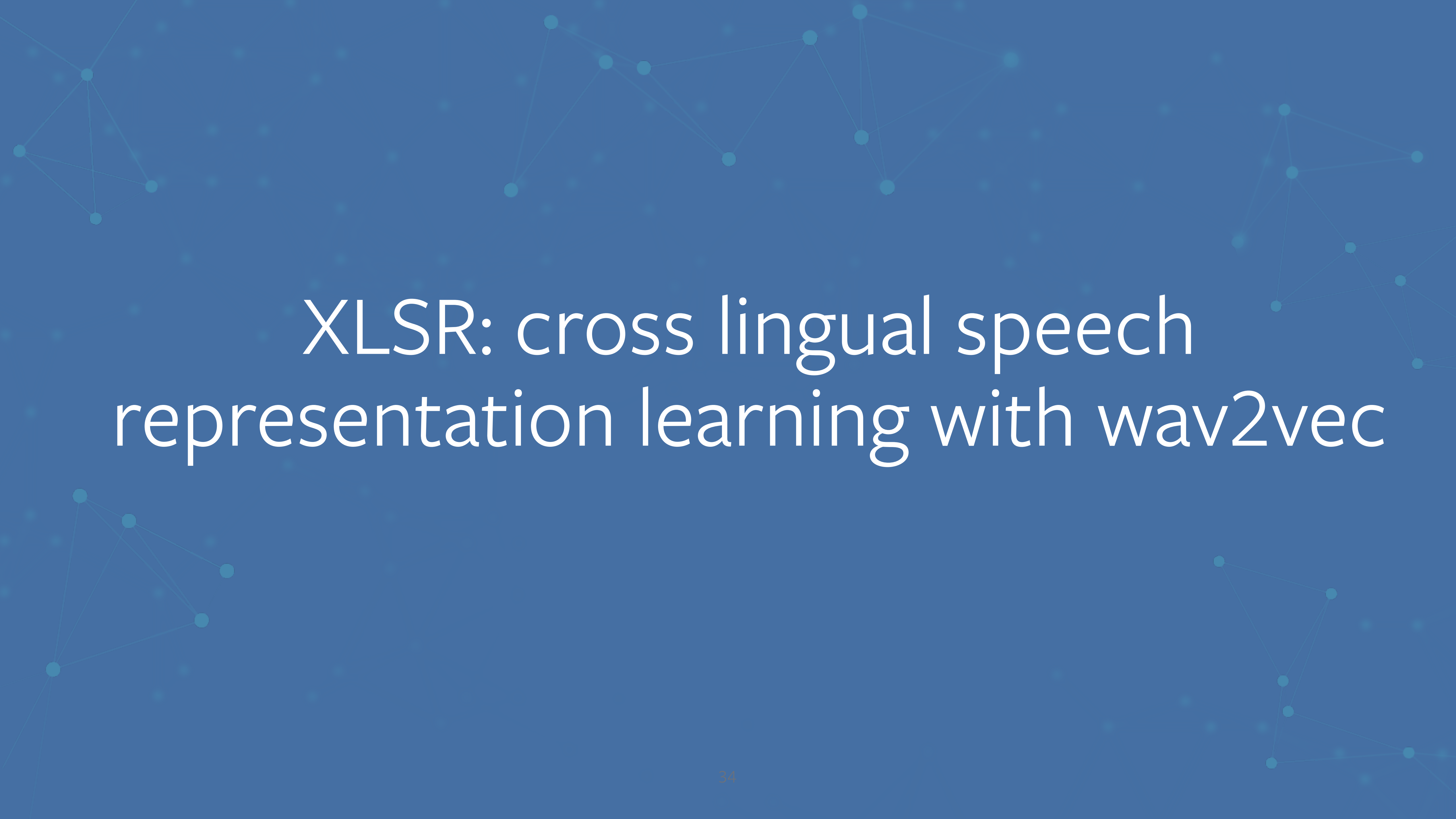
Amount of
labeled
data used



960h labeled

↑
10min
labeled



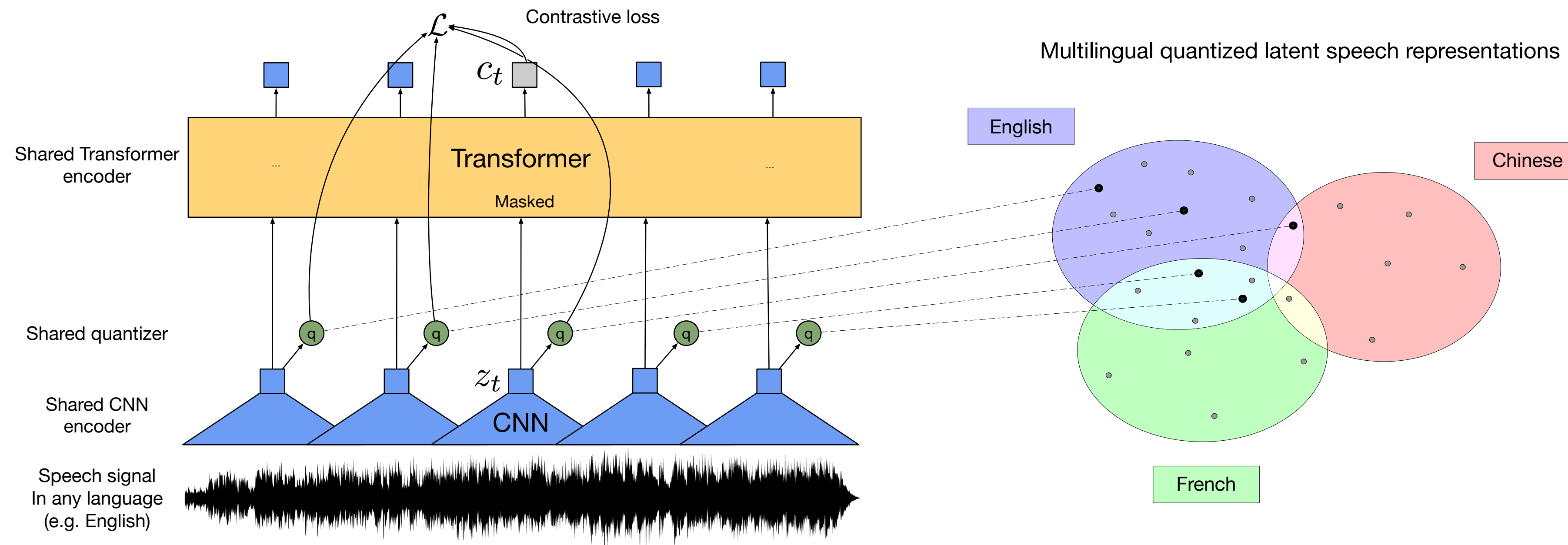
The background of the slide is a solid blue color with a faint, light blue network graph pattern. The graph consists of several clusters of nodes connected by thin lines, resembling a social network or a complex data structure. The nodes are small circles, and the lines are thin and light blue.

XLSR: cross lingual speech representation learning with wav2vec

Why *cross-lingual* self-supervised learning

- Little labeled data -> little unlabeled data
- Leverage unlabeled data from high-resource languages
- To improve performance on low-resource languages
- One model for each of the 6500 languages, for each domain? No.
- Instead: one pertained model for all languages

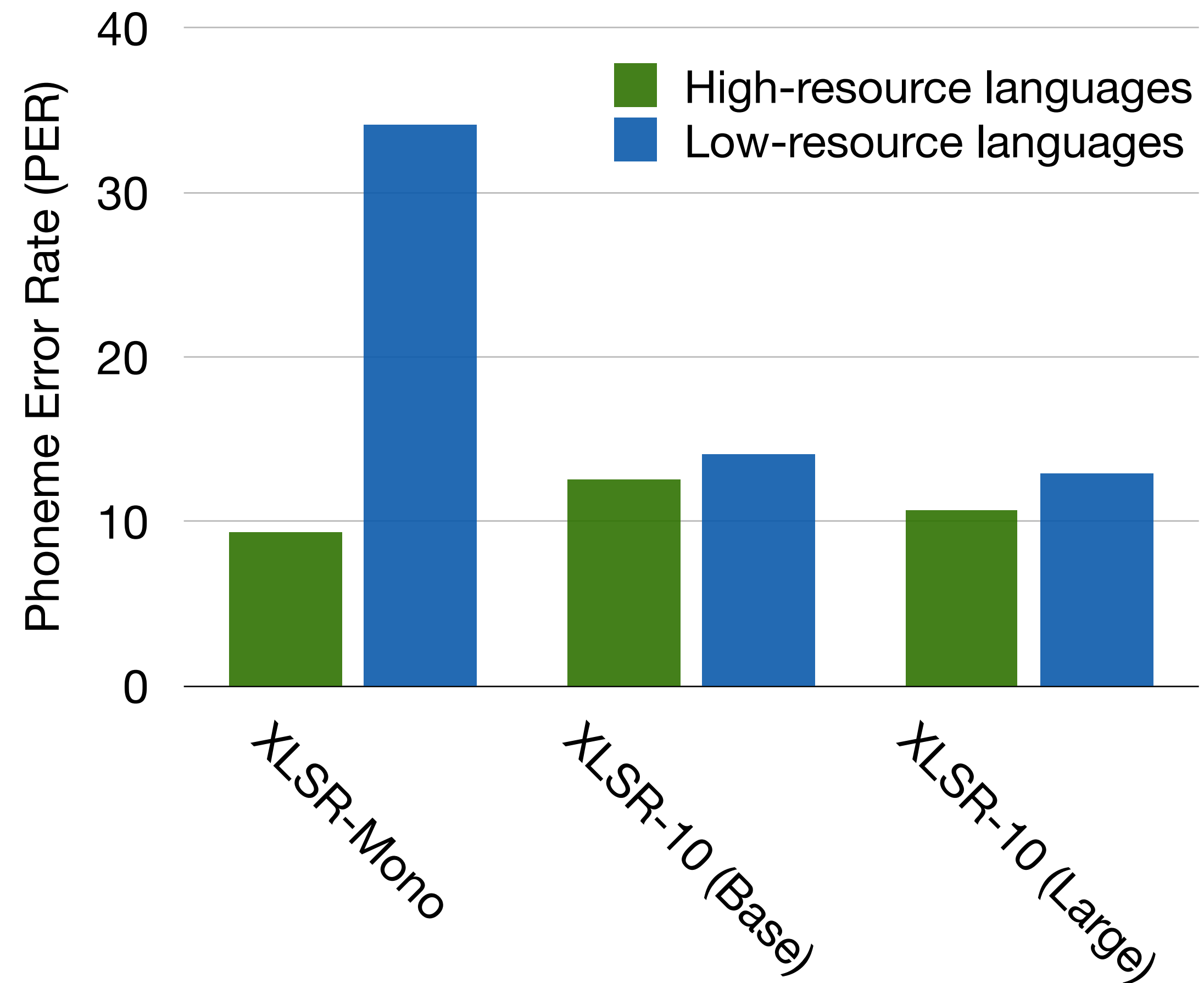
XLSR: cross lingual speech representation learning with wav2vec



XLSR: Results - cross-lingual transfer

Cross-lingual transfer = Train data from high-resource languages benefits low-resource languages.

CommonVoice results:



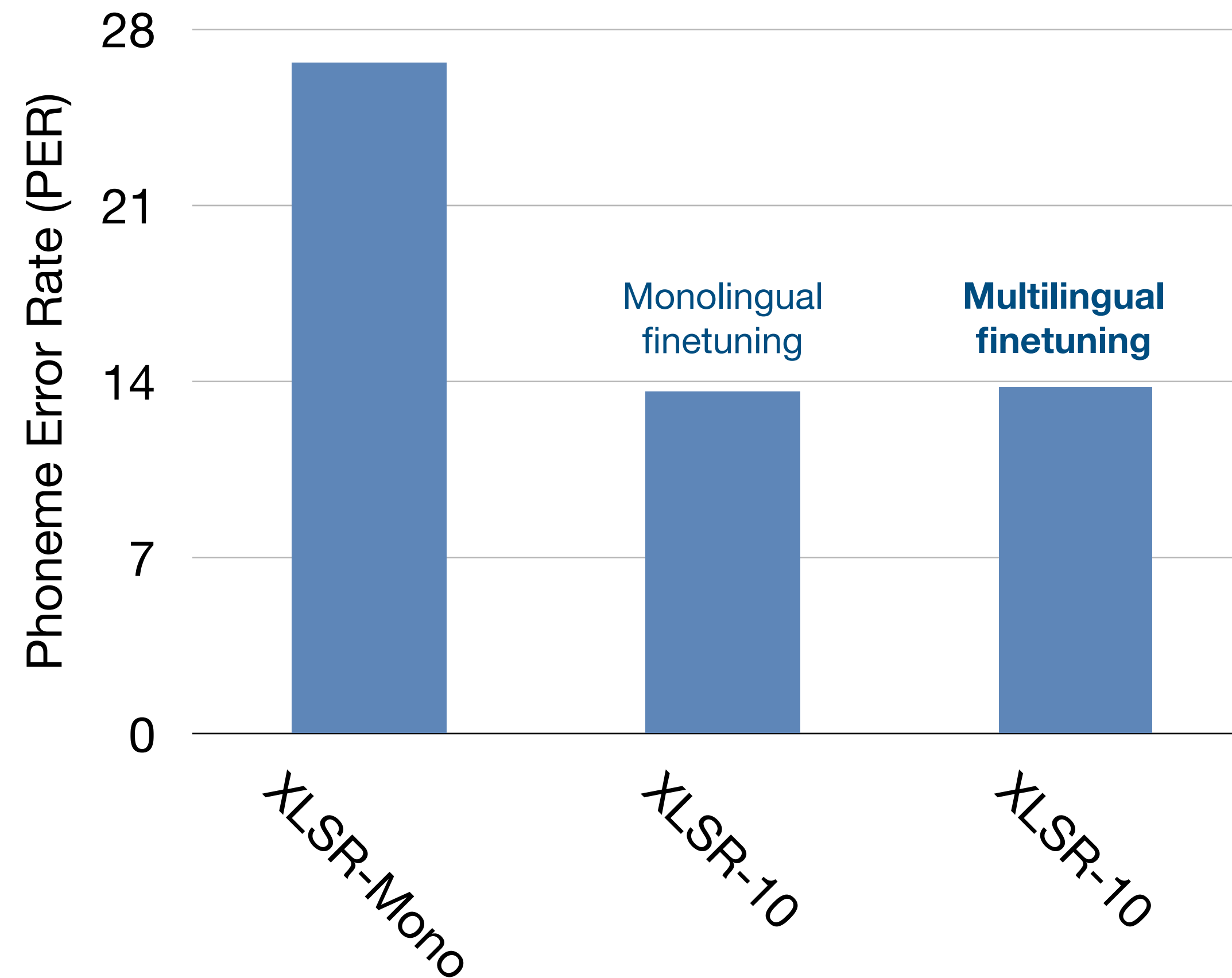
XLSR: Results - multilingual fine-tuning

Multilingual finetuning leads to *one model for all languages* with little loss in performance

XLSR: Results - multilingual fine-tuning

Multilingual finetuning leads to *one model for all languages* with little loss in performance

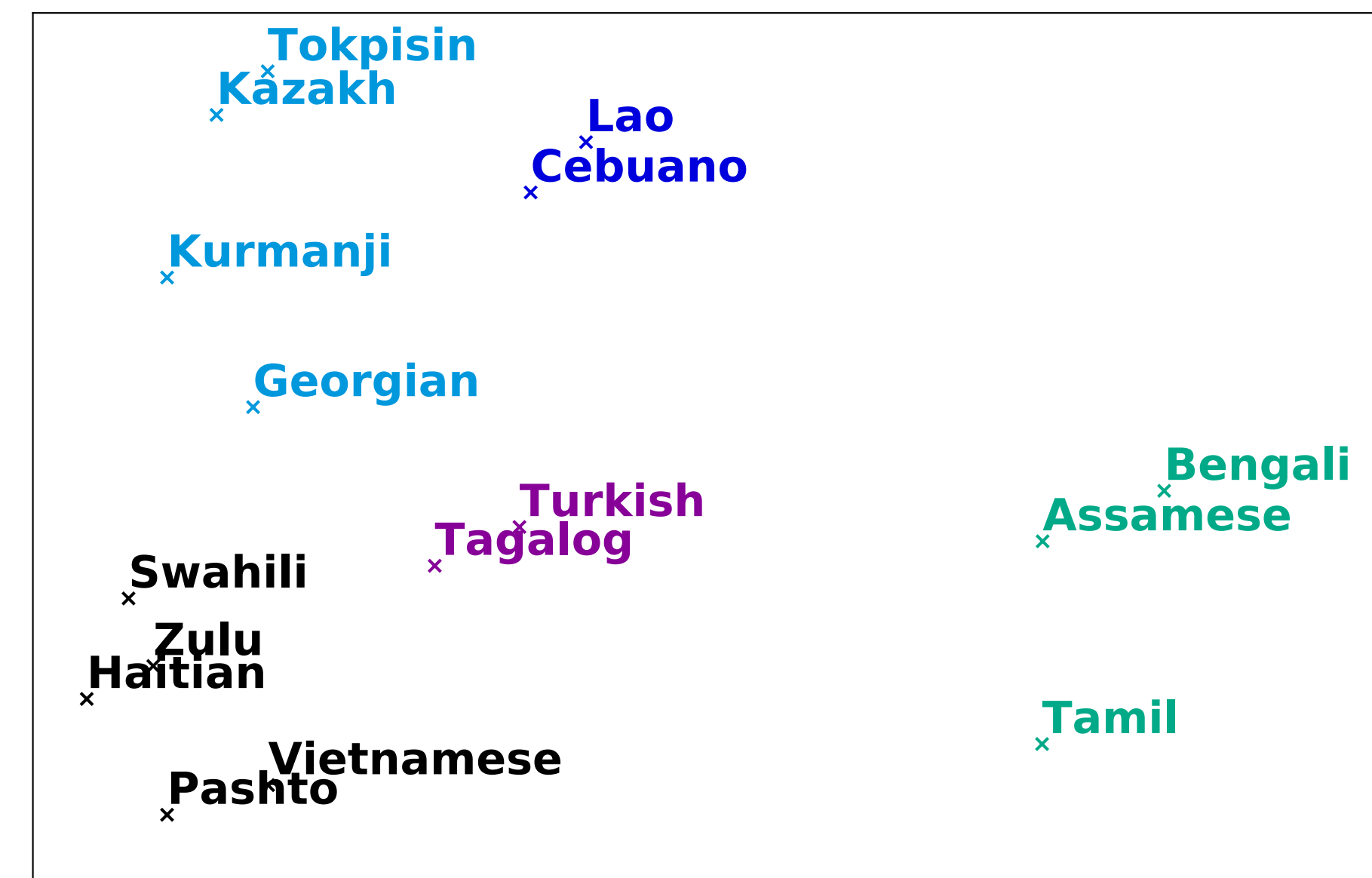
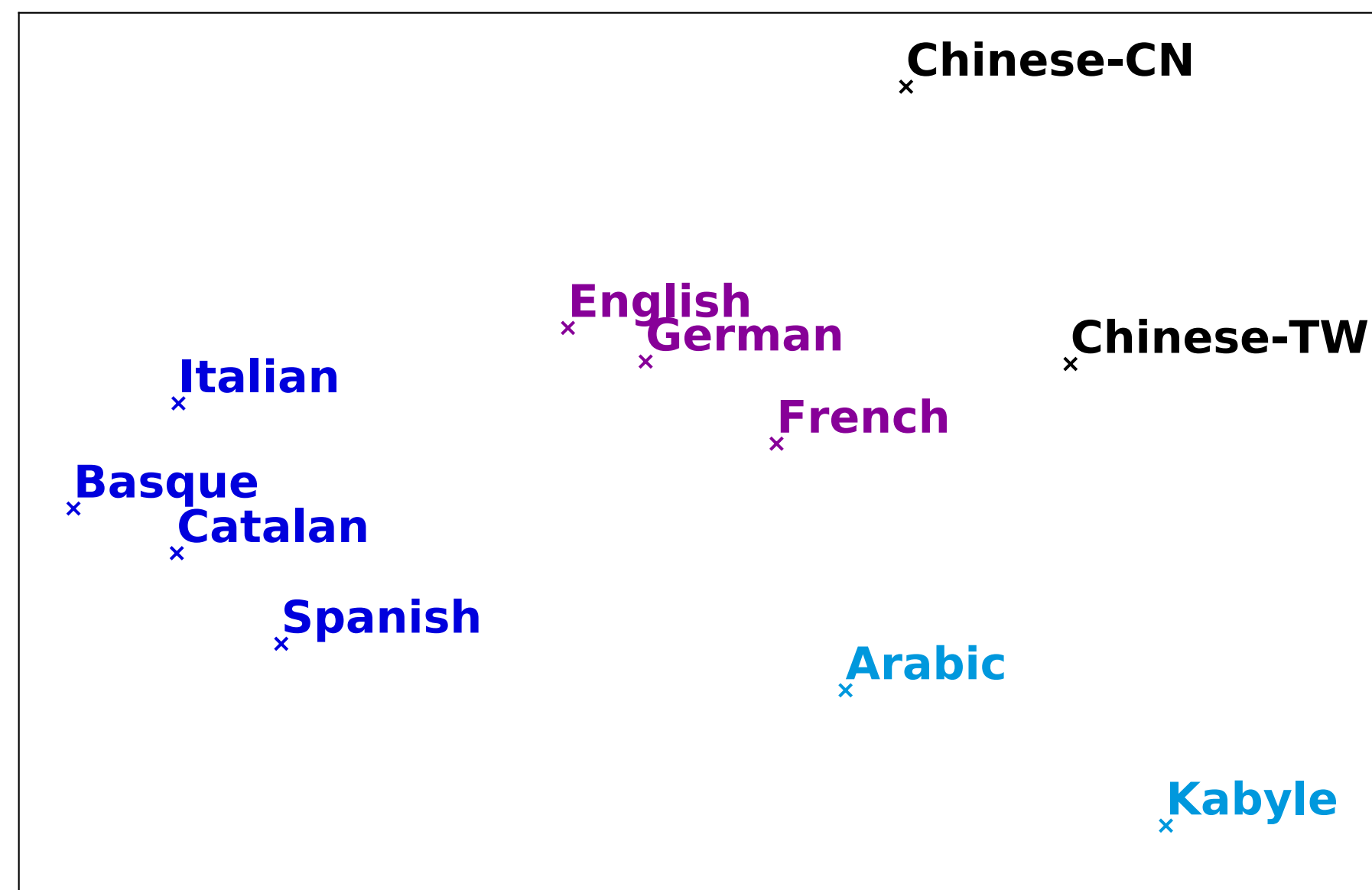
CommonVoice results:



XLSR: Analysis of discrete latent speech representations

PCA visualization of latent discrete representations from the multilingual codebook

Similar languages tend to share discrete tokens and thus cluster together



Unsupervised Speech Recognition

Unsupervised speech recognition

- Entirely remove need for labeled data
- Unsupervised machine translation works*, what about speech?
- Key problem: what are the units in the speech audio?

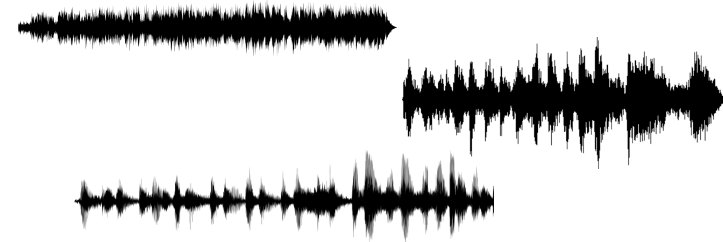


wav2vec Unsupervised: Key ideas

- Learn good representations of speech audio
- Unsupervised segmentation of the speech audio into phonemic units
- Learn mapping between speech segments and phonemes using adversarial learning

wav2vec Unsupervised

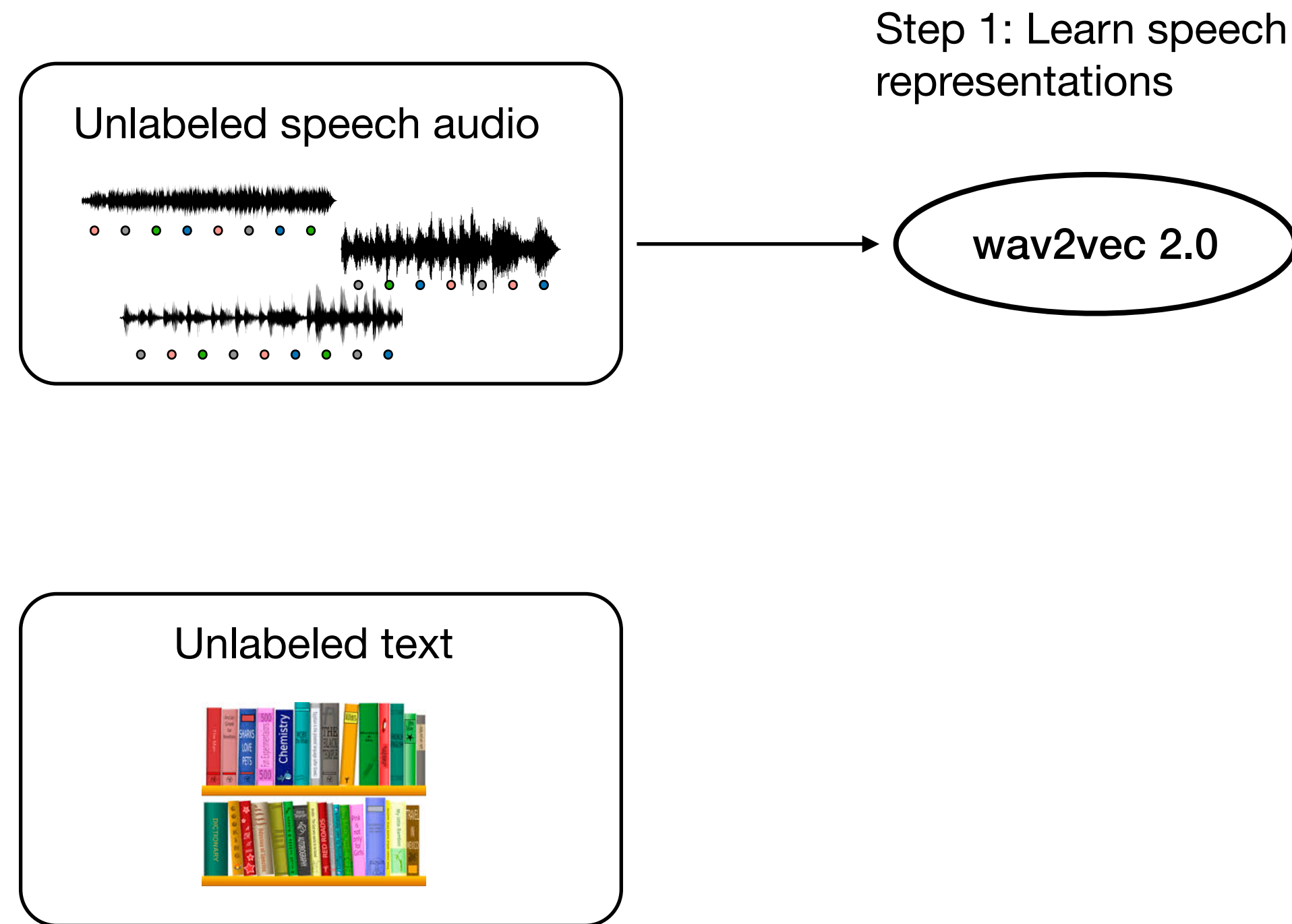
Unlabeled speech audio



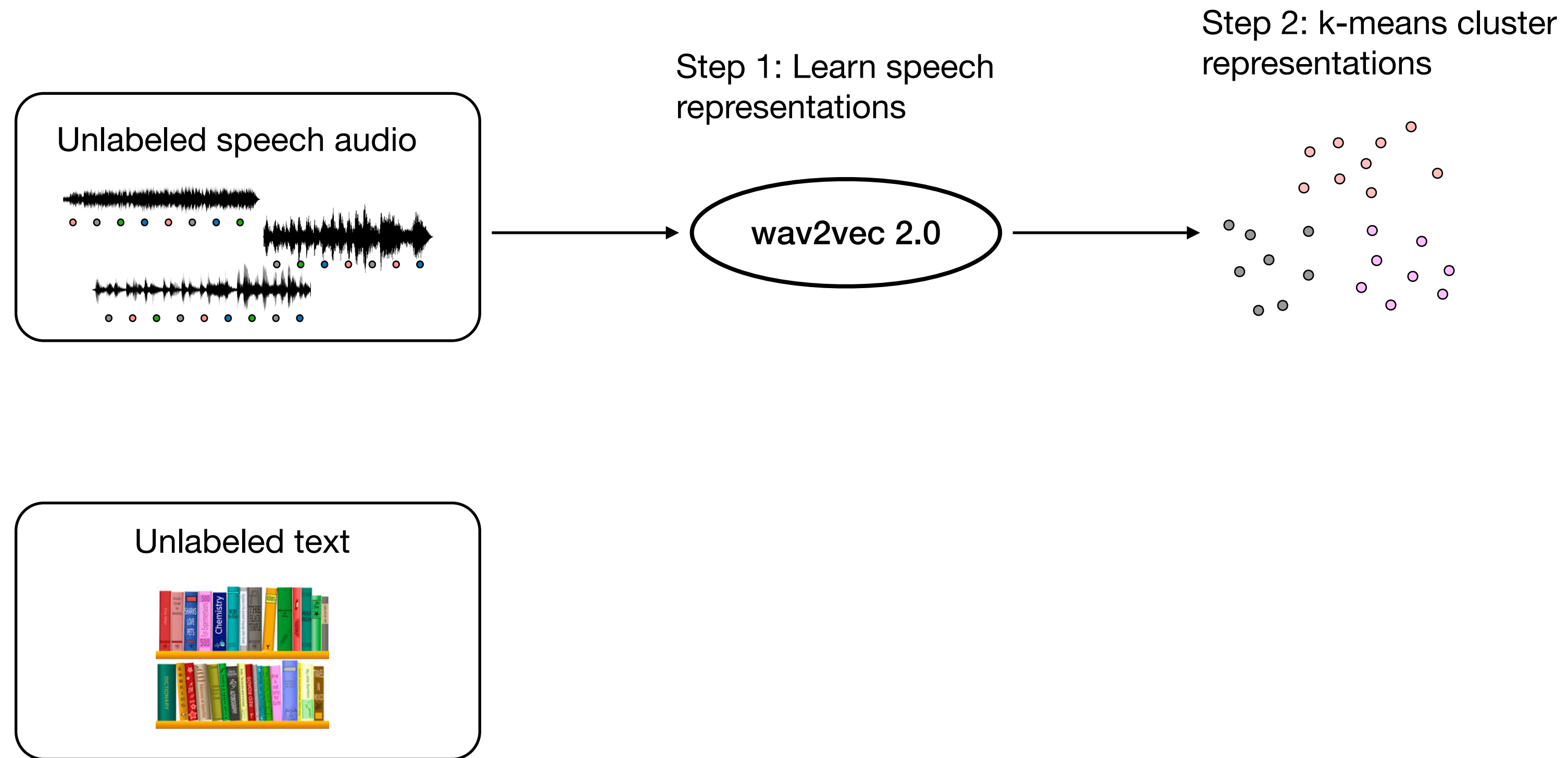
Unlabeled text



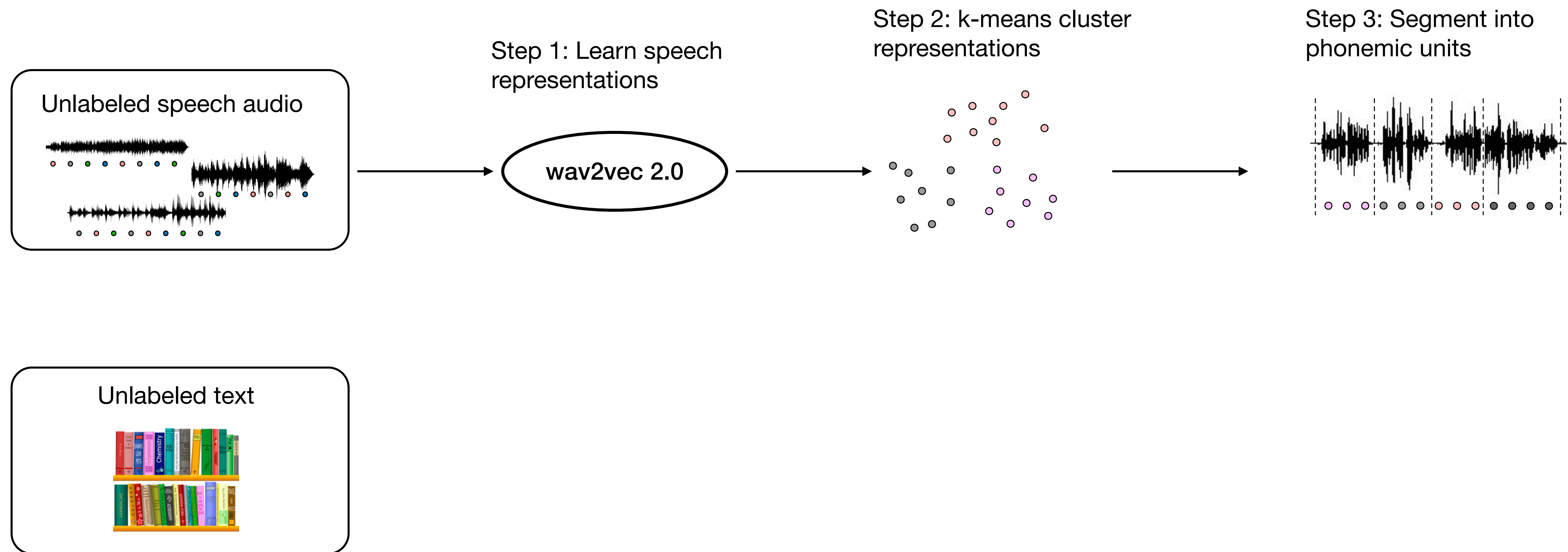
wav2vec Unsupervised



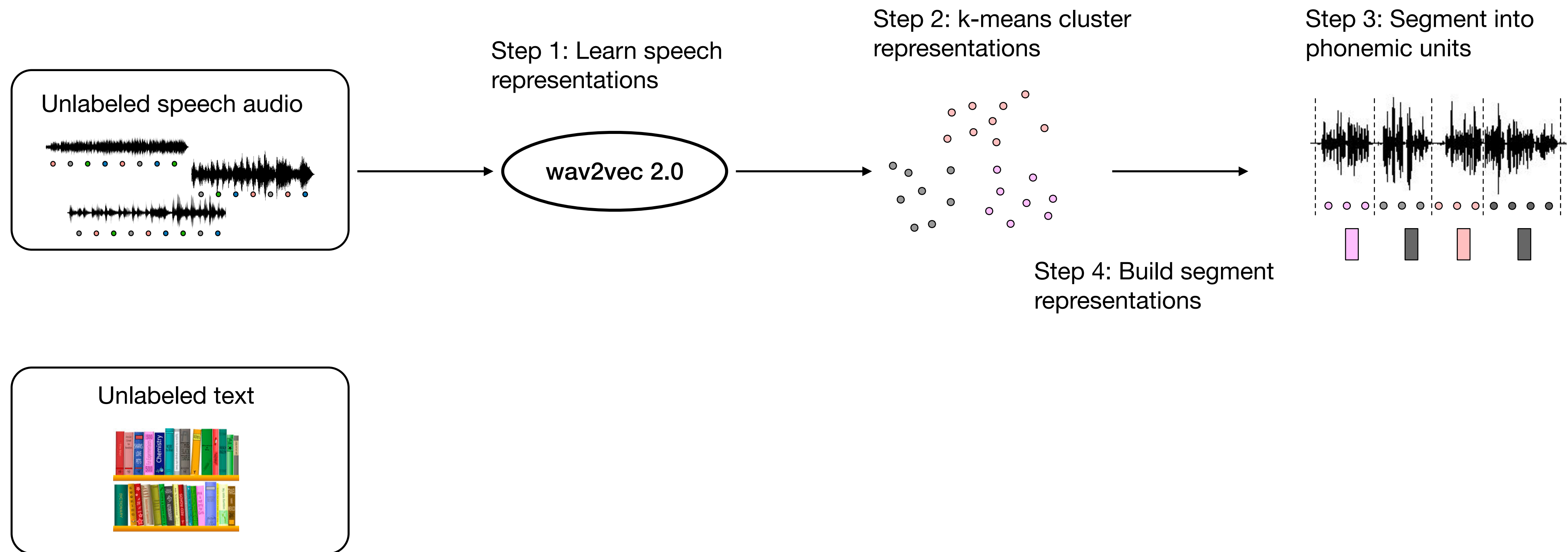
wav2vec Unsupervised



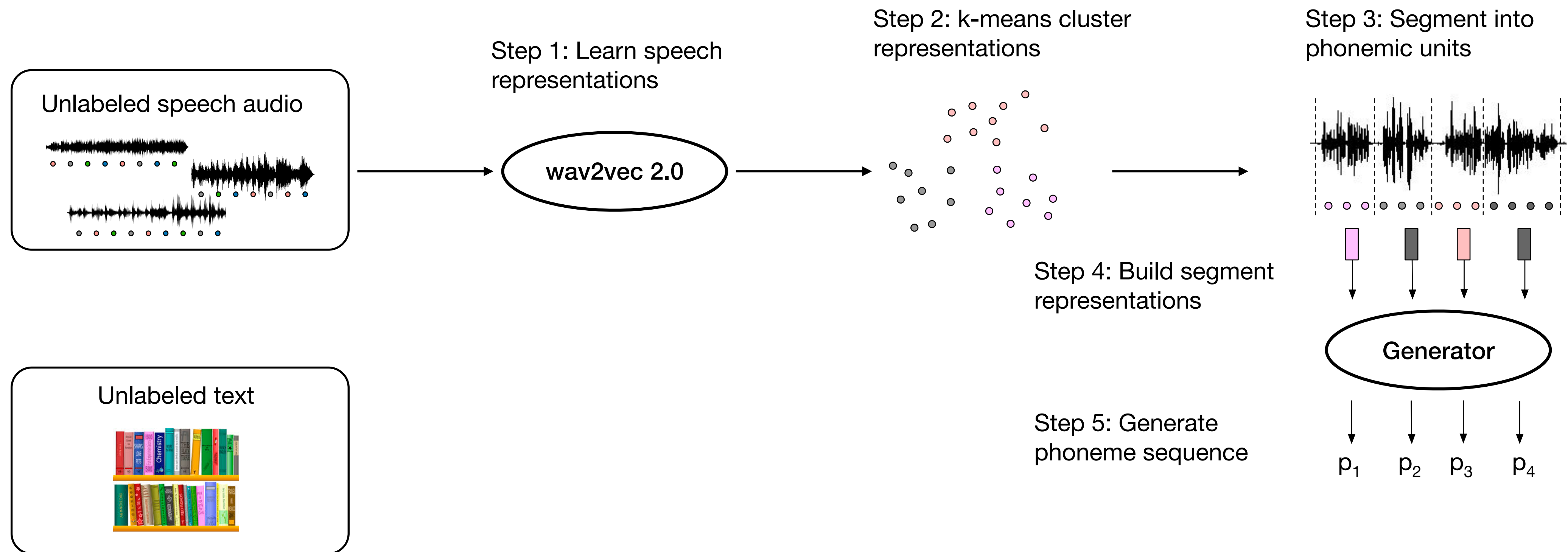
wav2vec Unsupervised



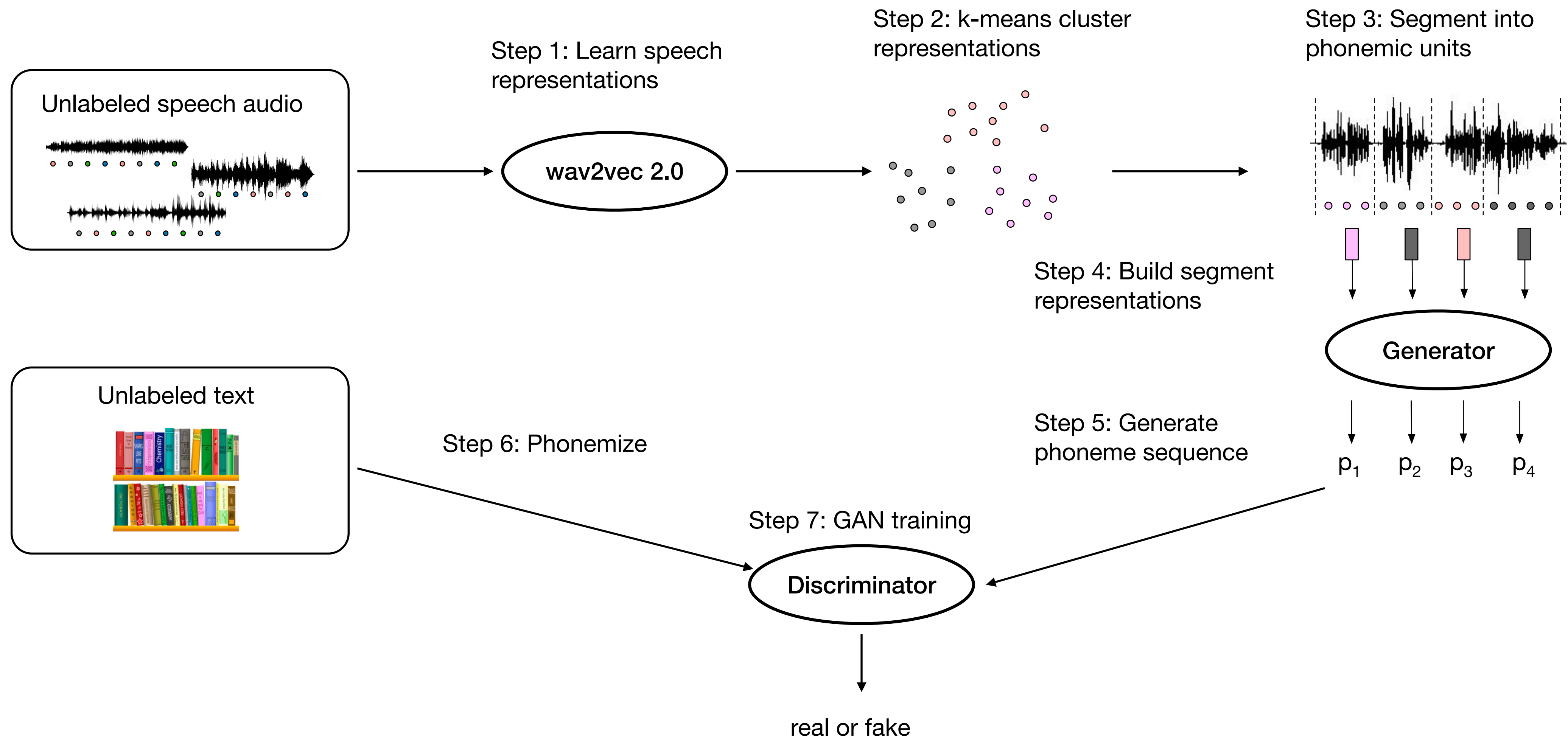
wav2vec Unsupervised



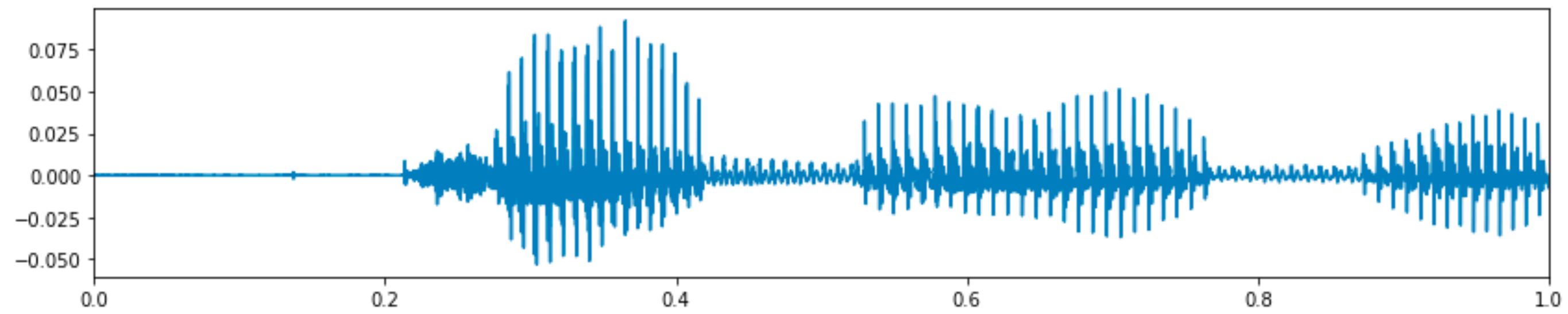
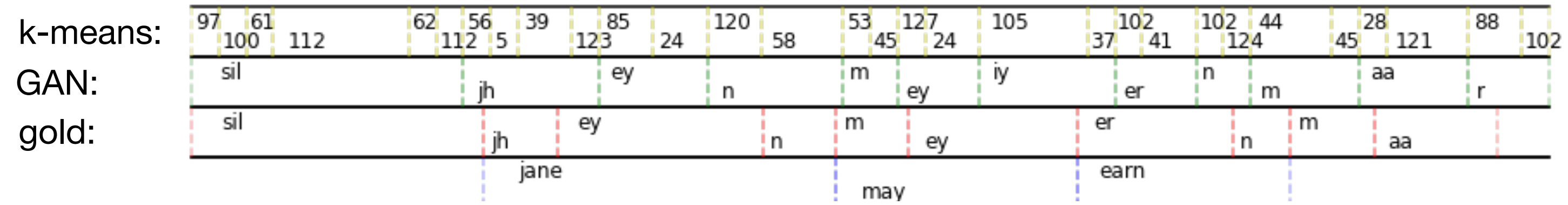
wav2vec Unsupervised



wav2vec Unsupervised



Simple segmentation



Text data pre-processing

Unlabeled text



he

spoke

soothingly

Text data pre-processing

Unlabeled text



he

spoke

soothingly

Phonemize

hh iy

s ow k

s uw dh ih ng l iy

Text data pre-processing

Unlabeled text



he

spoke

soothingly

Phonemize

sil

hh iy

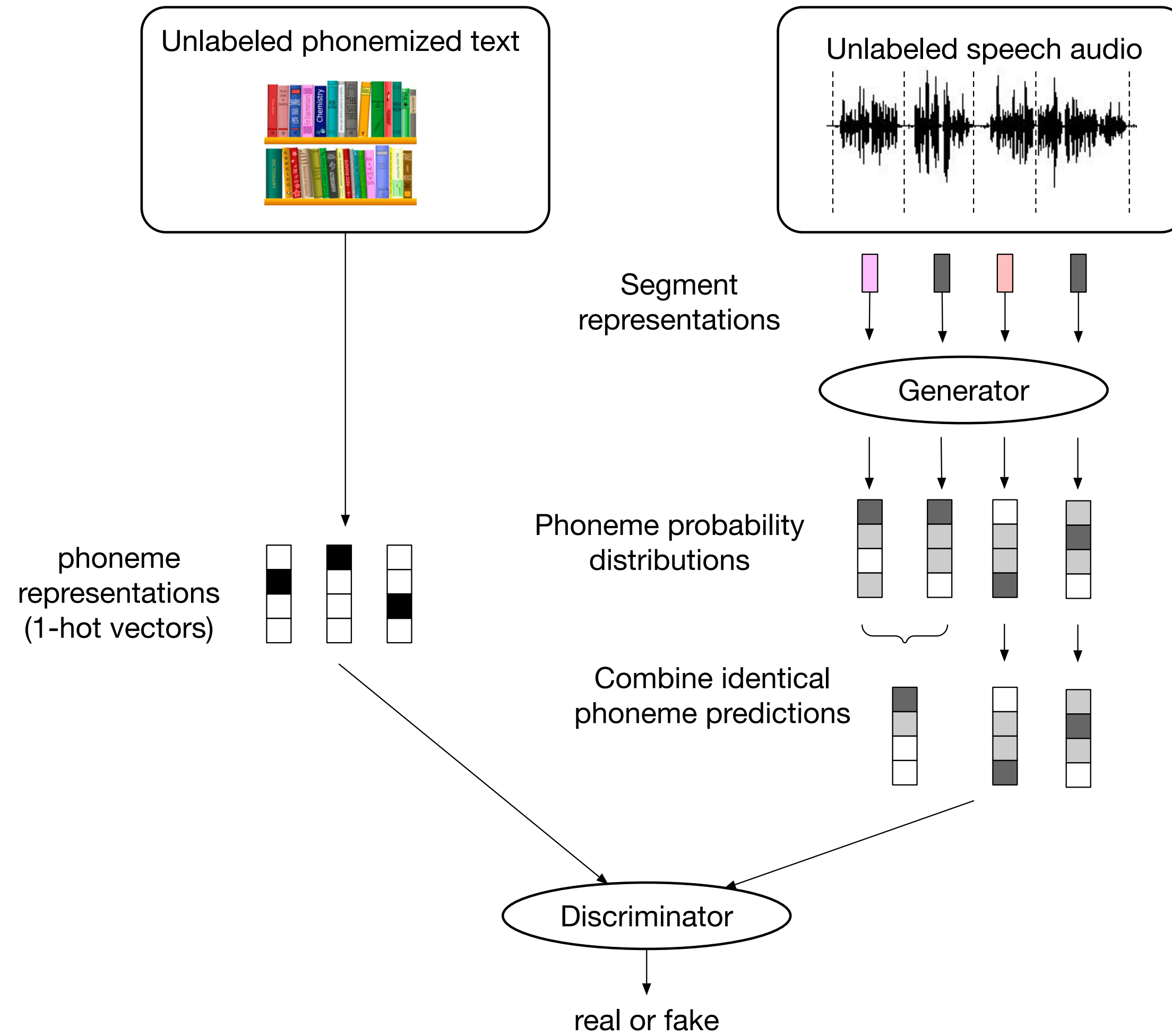
s ow k

s uw dh ih ng l iy

sil

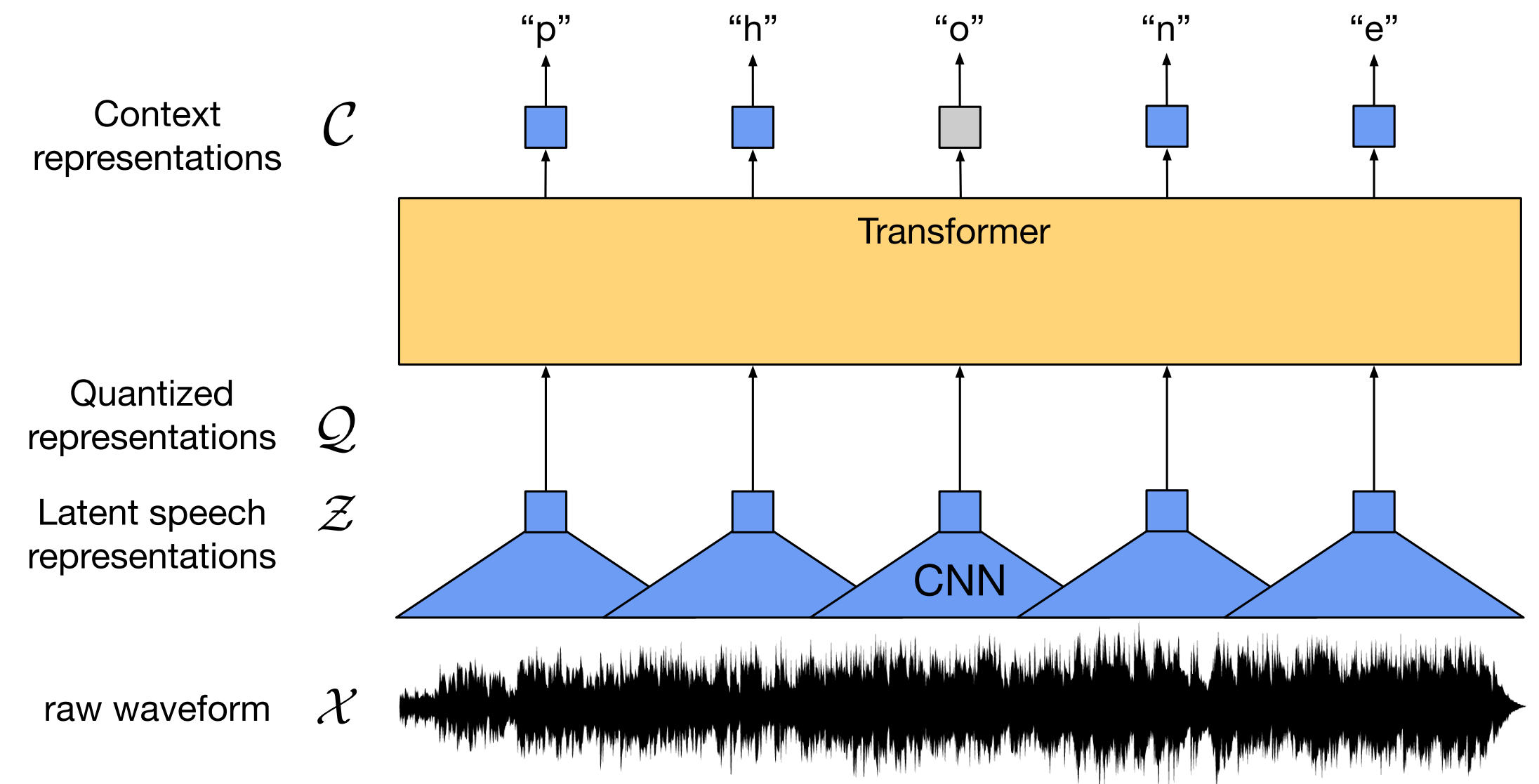
Silence insertion

GAN inputs



Generator / Discriminator

- Generator: 1 layer CNN with 90k parameters
w2v features frozen
- Discriminator: 3 layer CNN
- Train time: 12-15h on a single V100

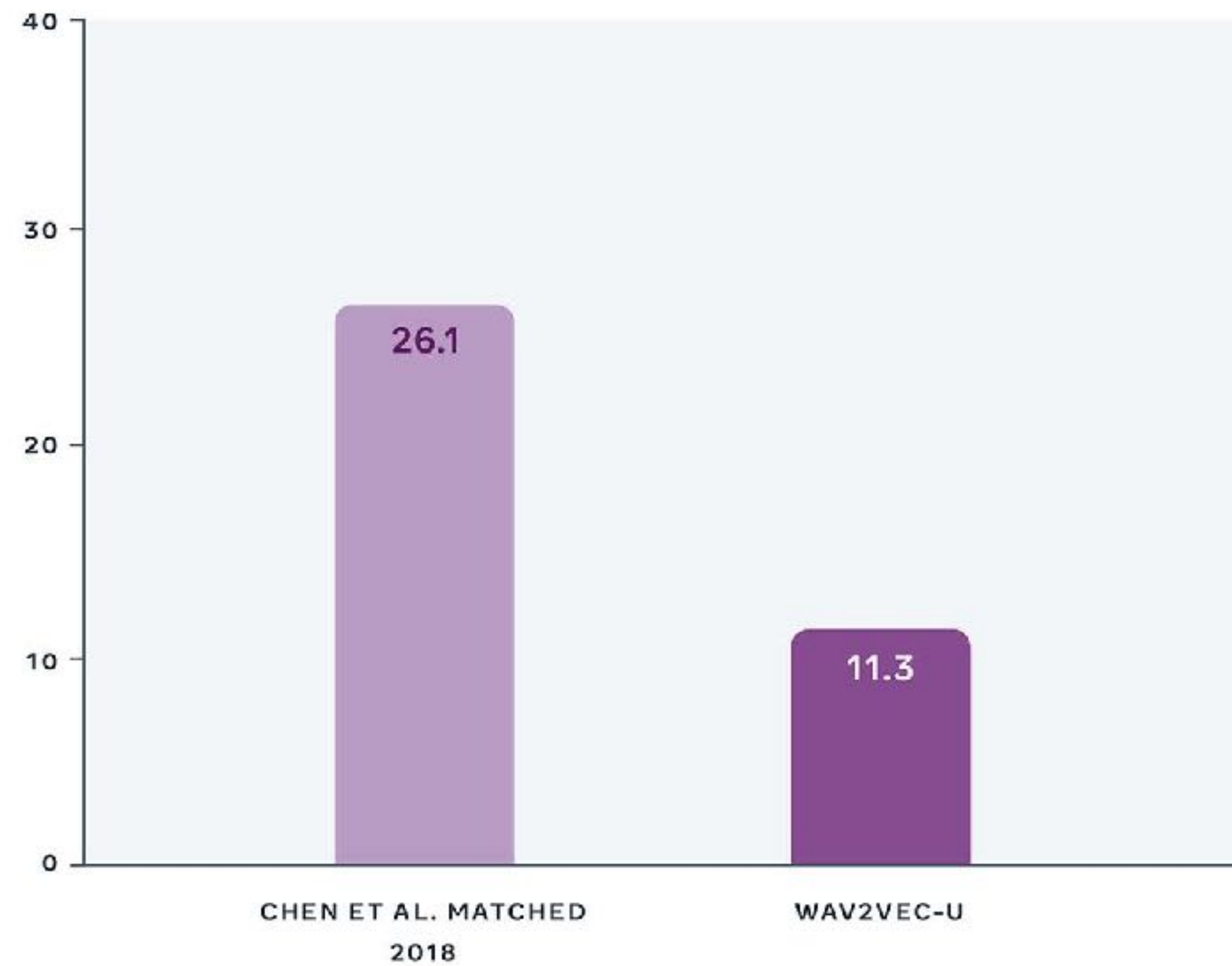


Training details

- Unsupervised metric for early stopping, hyper-parameter selection
- Self-training after GAN training (HMM and fine-tuning w2v)

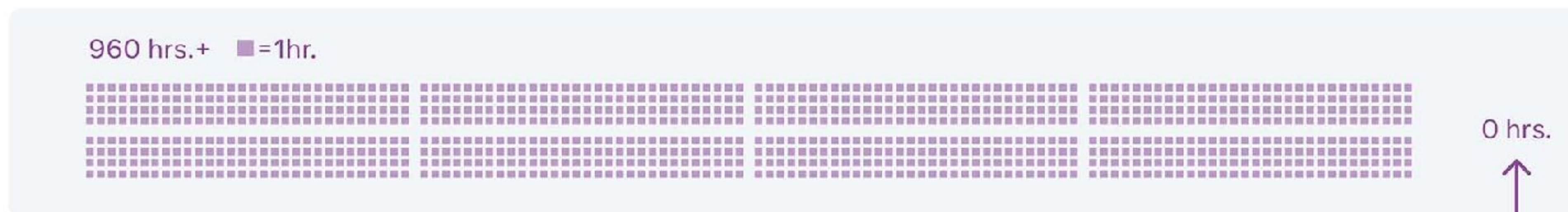
Comparison to prior unsupervised work

Phoneme error rate

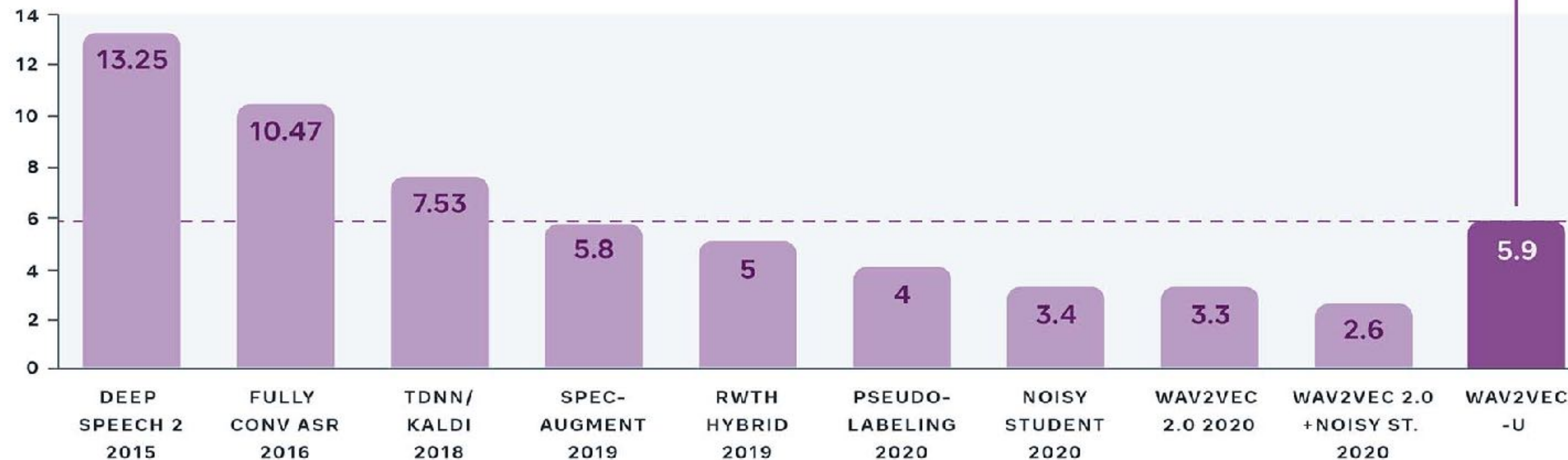


Comparison to best supervised systems

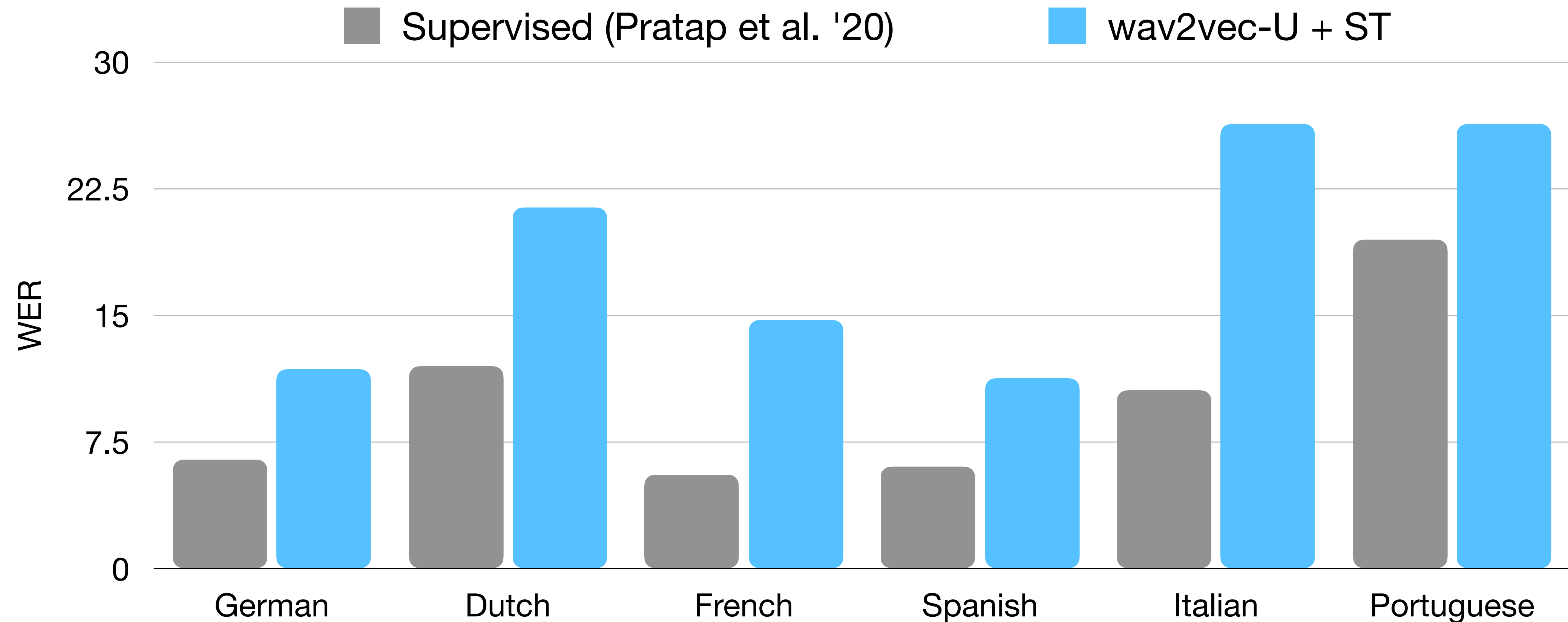
Amount of labeled data used



Word error rate

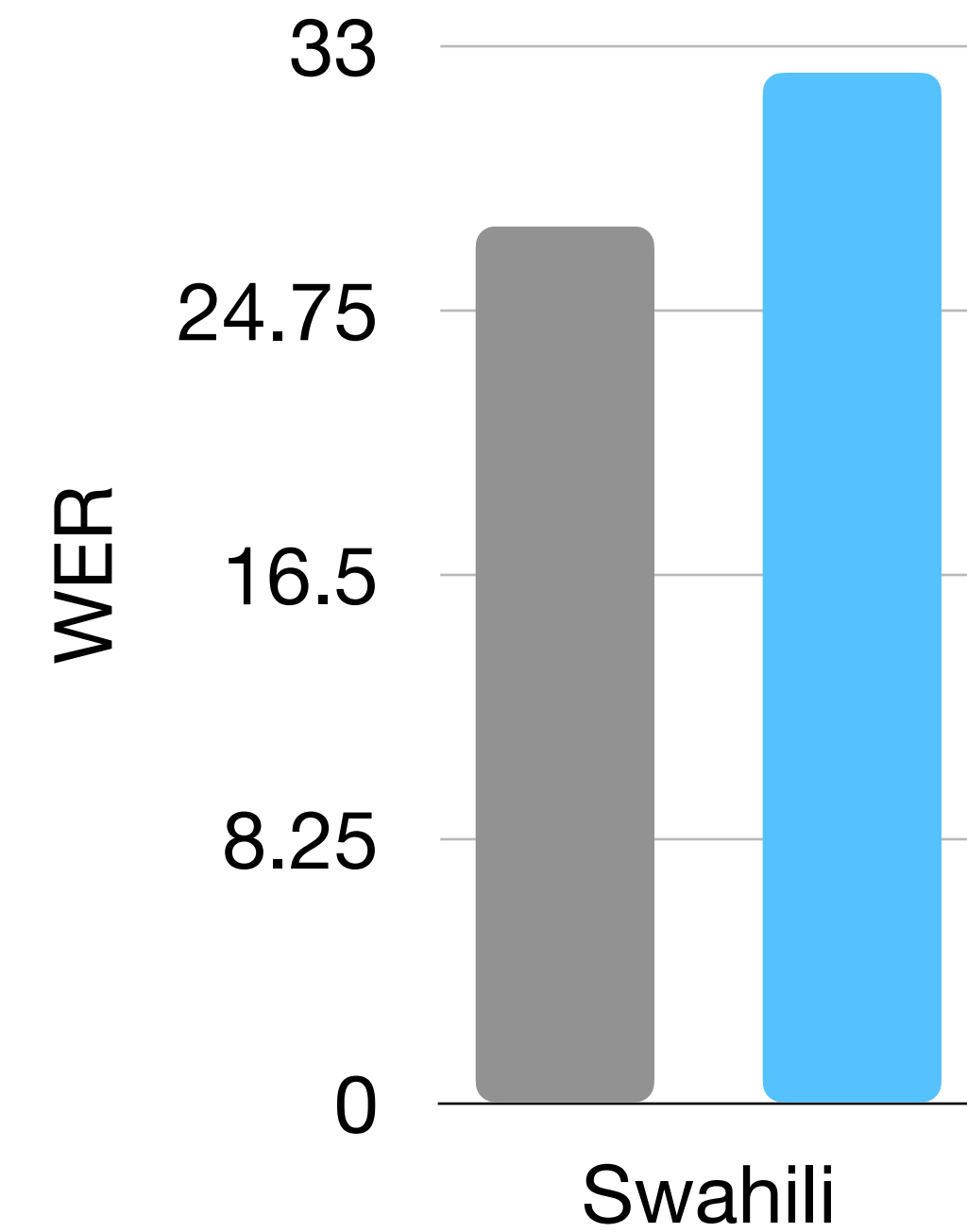
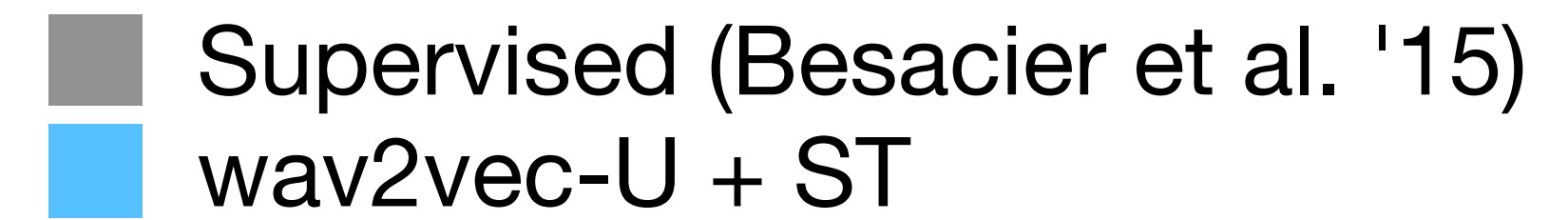
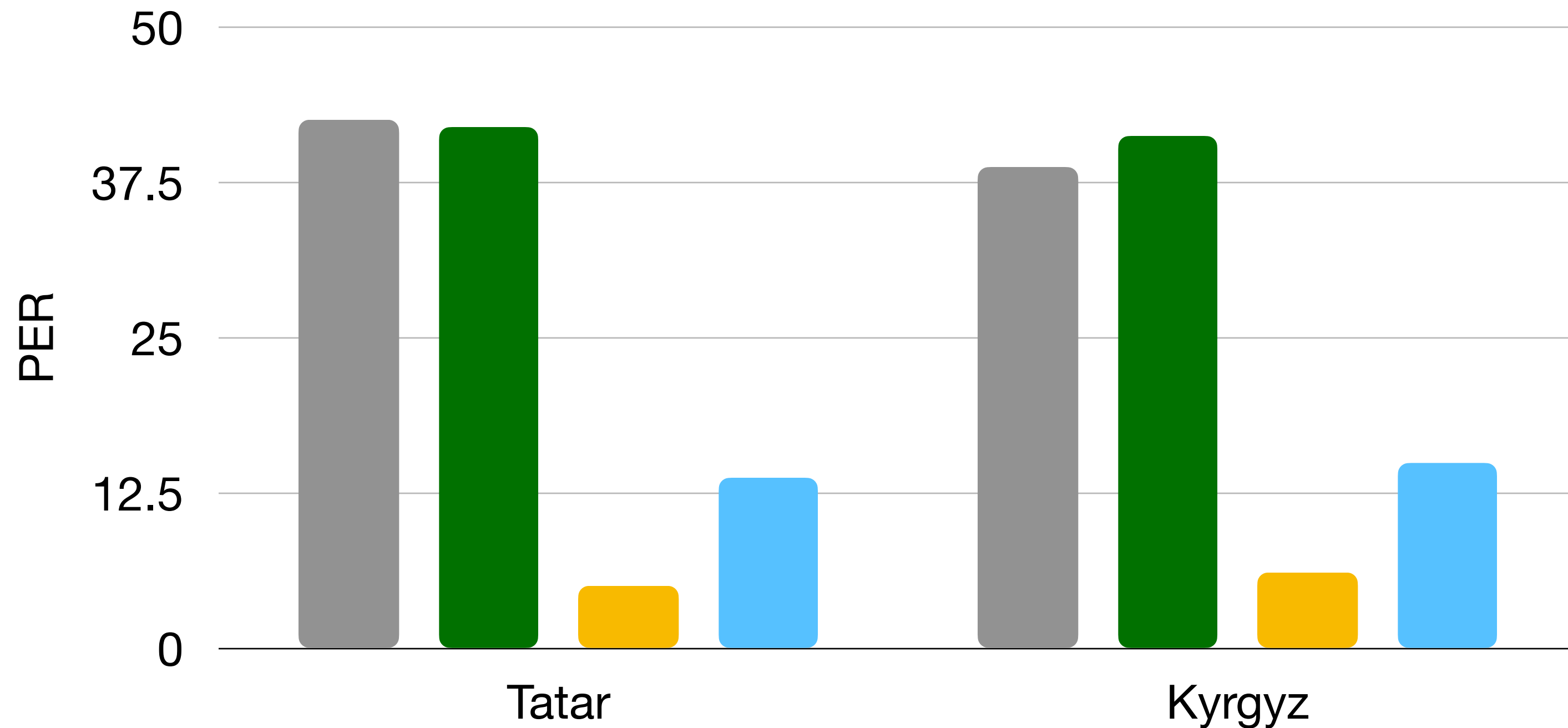


Other languages



MLS benchmark, wav2vec-U used only 100h of unlabeled data but there is up to 2k hours for some languages.

Low-resource languages



*wav2vec-U uses much less speech audio than prior work:
1.8h vs. 17h for Kyrgyz, 4.6h vs. 17h for Tatar

Discussion

- Very lightweight approach (except for wav2vec 2.0)
- Why does it work? Good audio features are main driver of performance
- Phonemizer still required
- Segment construction

Conclusion

- Pre-training for speech works very well in both low-resource and high-resource setup.
- Cross-lingual training improves low-resource languages.
- Enable speech models with very little or even no labeled training data
- Make speech technology more ubiquitous and robust
- Code and models are available in the fairseq GitHub repo + Hugging Face.



Thank you



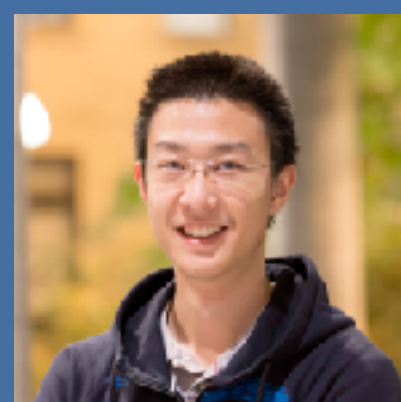
Alexei Baevski



Alexis Conneau



Steffen Schneider



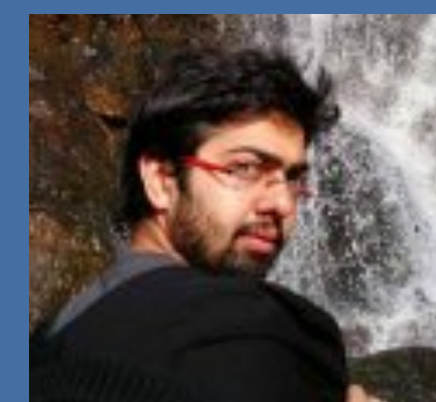
Henry Zhou



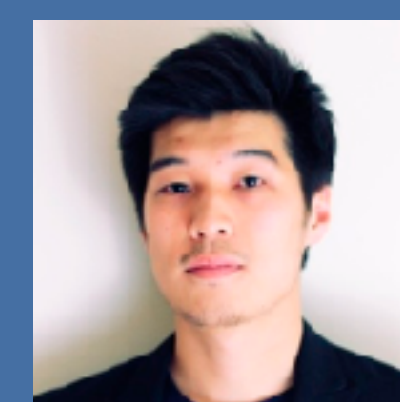
Abdelrahman
Mohamed



Anuroop
Sriram



Naman
Goyal



Wei-Ning Hsu



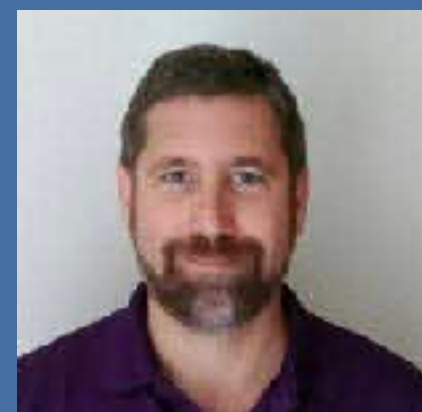
Michael Auli



Kritika Singh



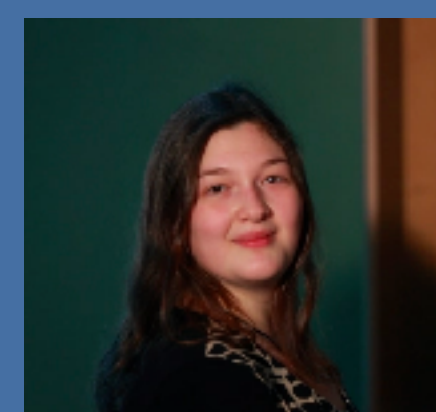
Yatharth Saraf



Geoffrey Zweig



Qiantong Xu



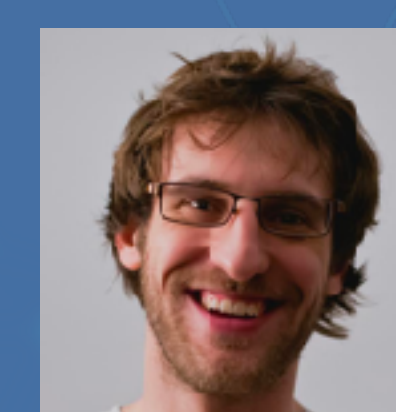
Tatiana
Likhomanenko



Paden
Tomasello



Ronan
Collobert



Gabriel
Synnaeve