# CS11-737 Multilingual NLP
# Speech Pre-training

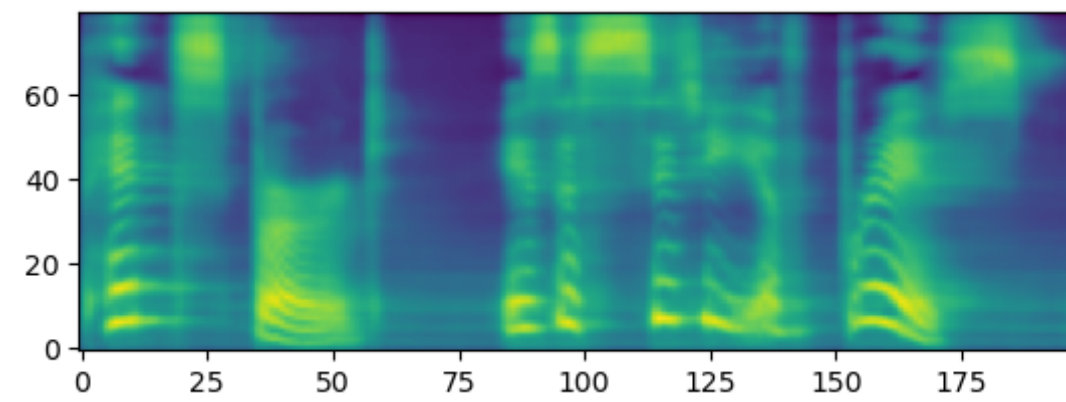Lei Li

https://lileicc.github.io/course/11737mnlp23fa/

**Carnegie Mellon University**
**Language Technologies Institute**

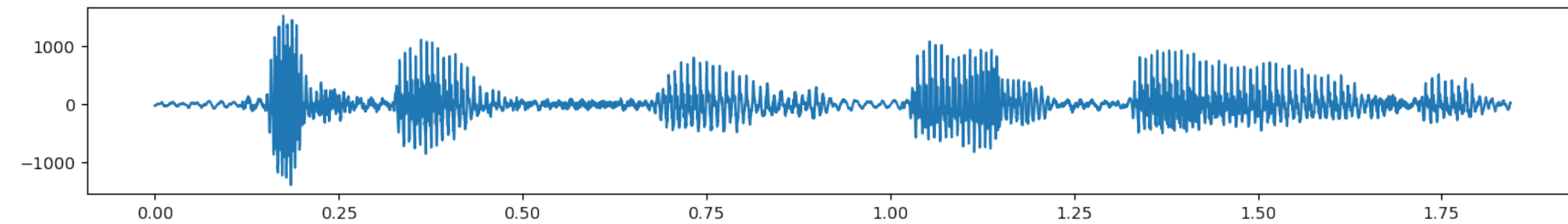# Feature Extraction for Speech Recognition

"Pittsburgh is a city of bridge"
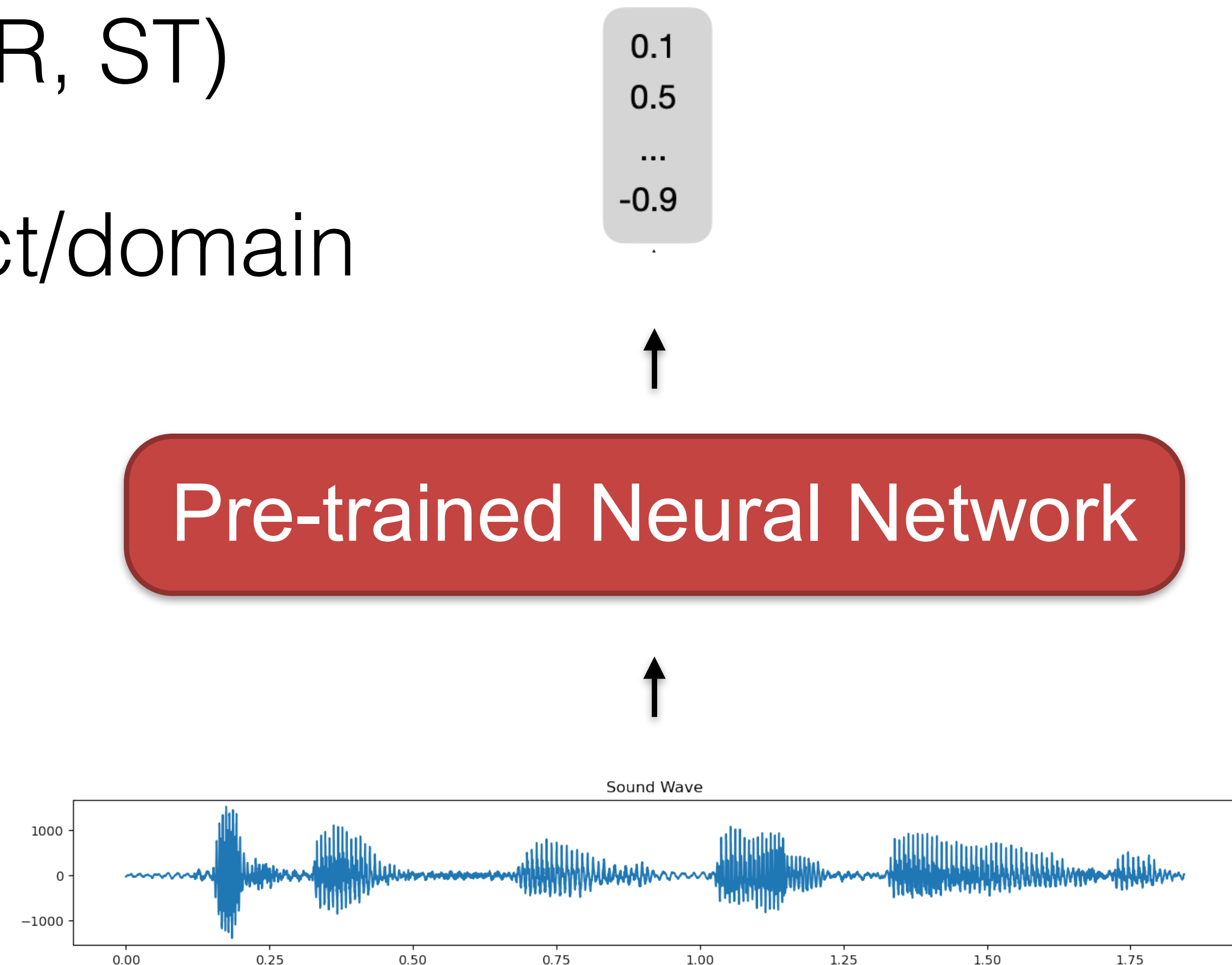
↑

**Neural Network**

↑

MFCC

Sound Wave

- need 1,000+ hours of transcribed data to train a good ASR system

- how to generalize to many languages/ dialects?

2

# Self-supervised Speech Representation Learning

- Self-supervised Training on unlabeled audio data

- generalize to many tasks (ASR, ST)

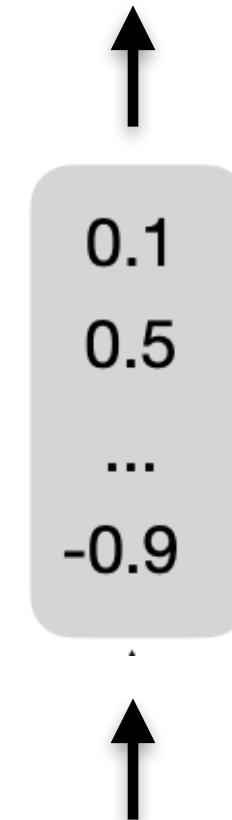- generalize to language/dialect/domain

0.1
0.5
...
-0.9

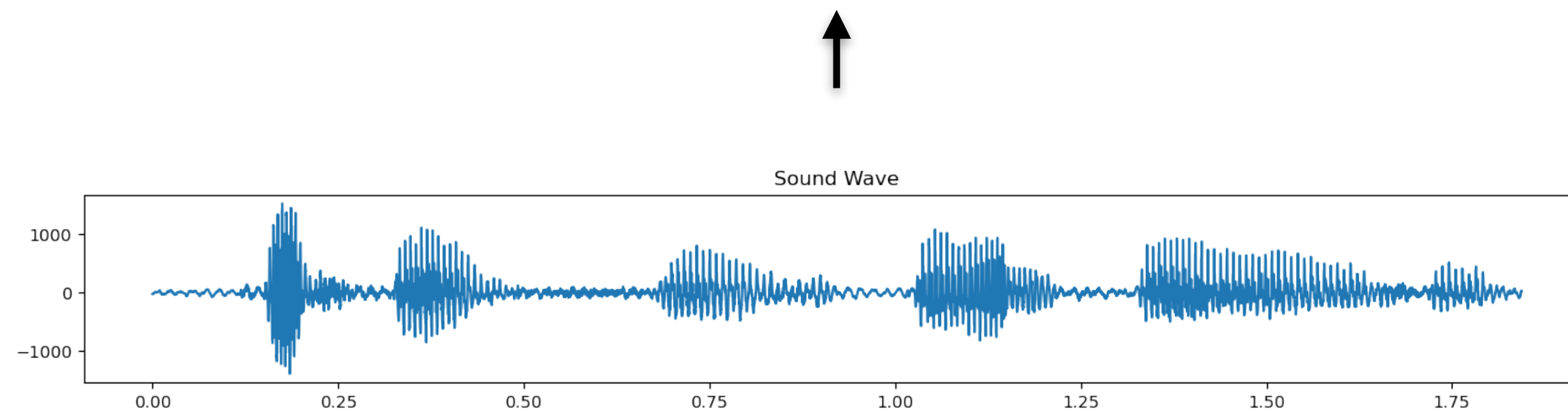**Pre-trained Neural Network**

Sound Wave

3

# Transfer to Downstream Tasks

Fine-tuning

Task-specific network

ASR
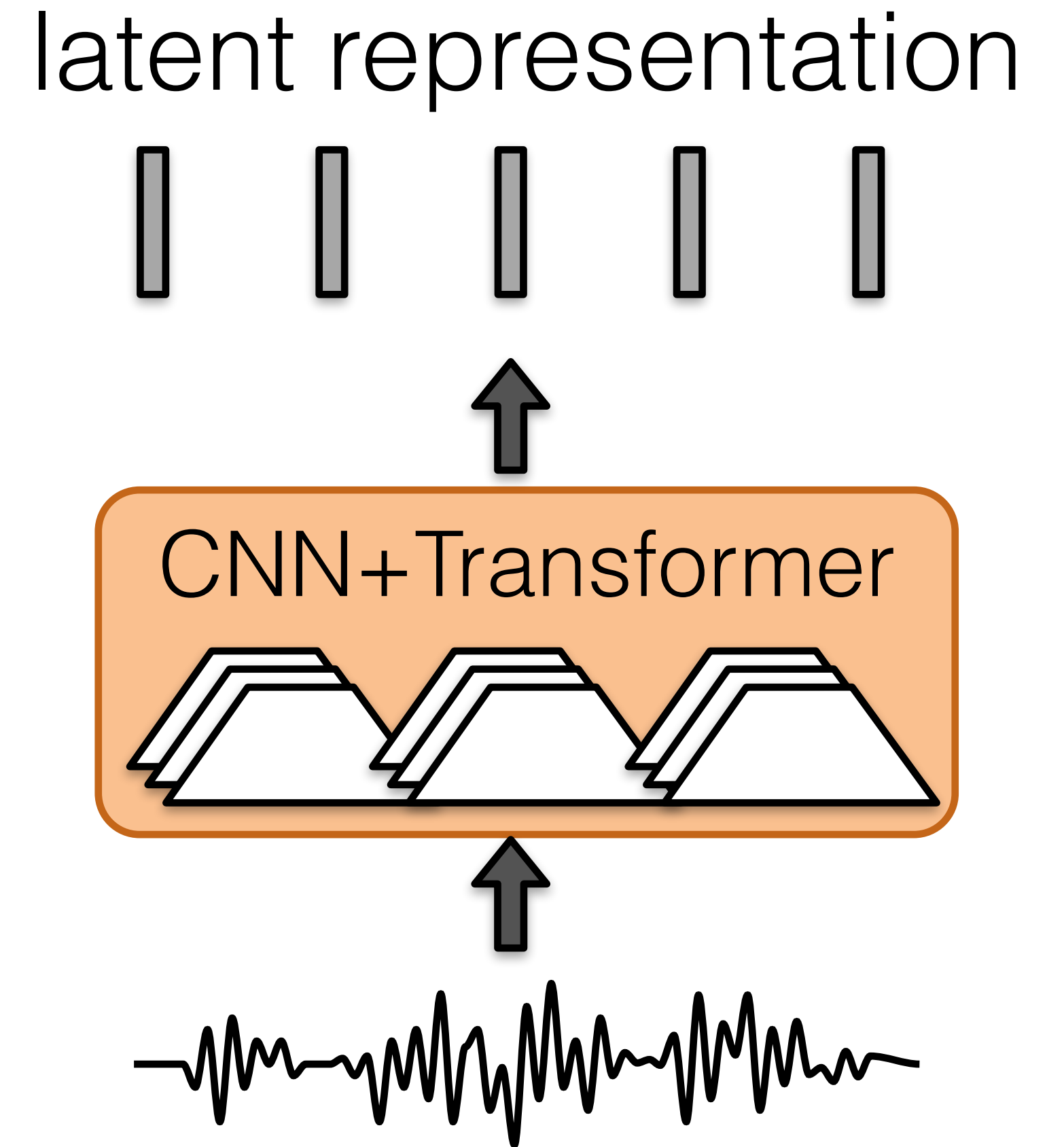Speech Translation

0.1
0.5
...
-0.9

Pre-trained Neural Network

Sound Wave
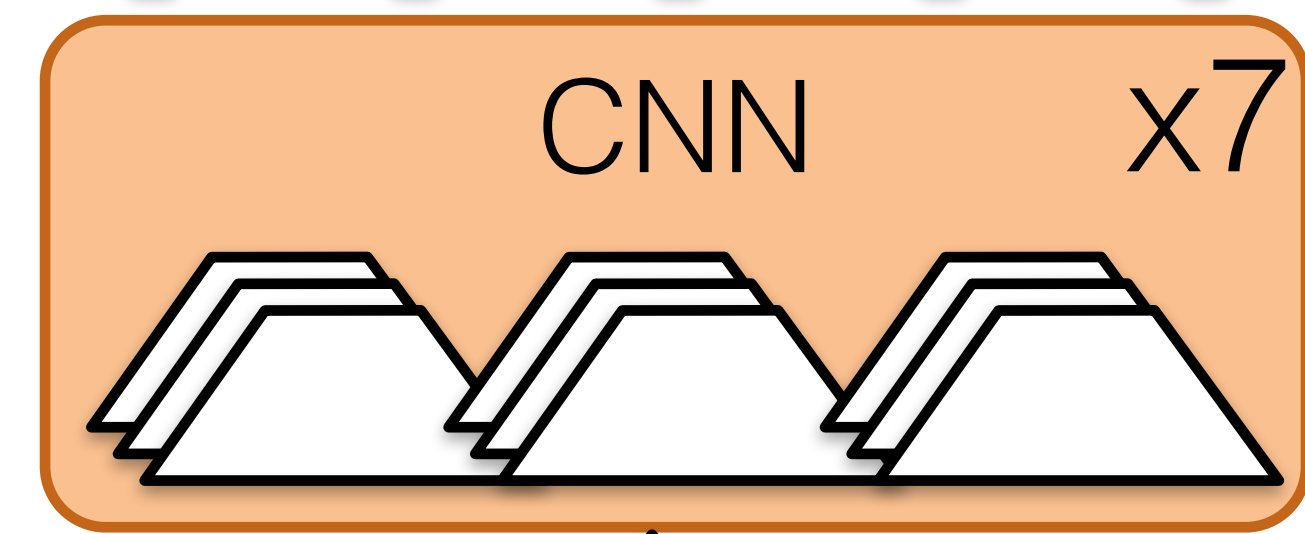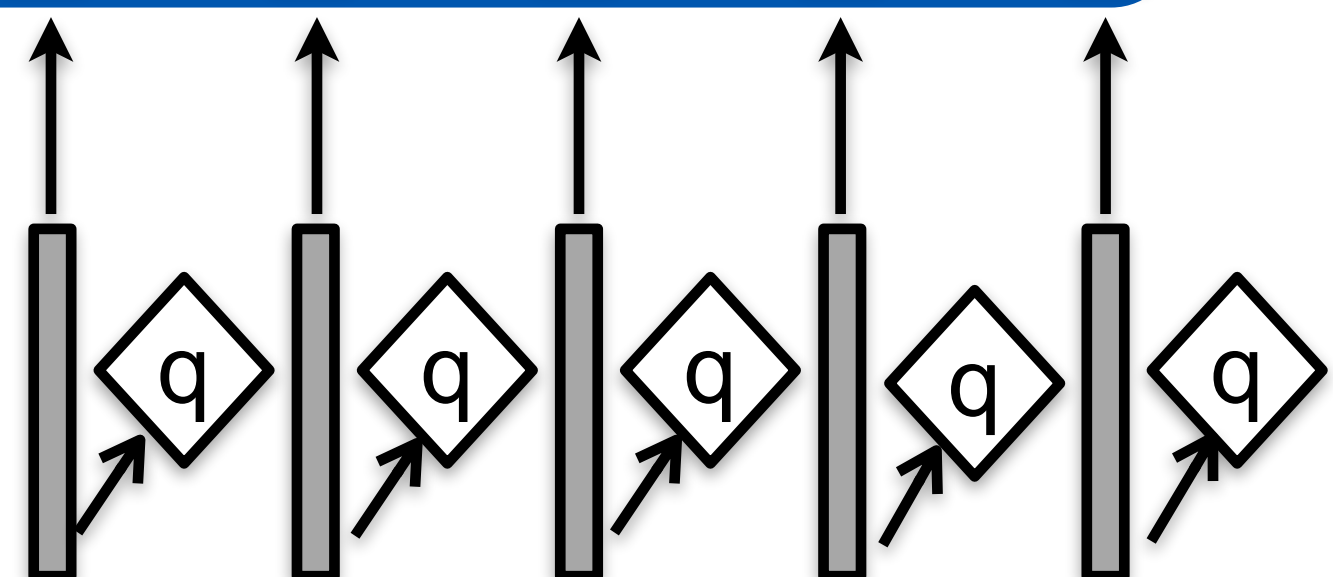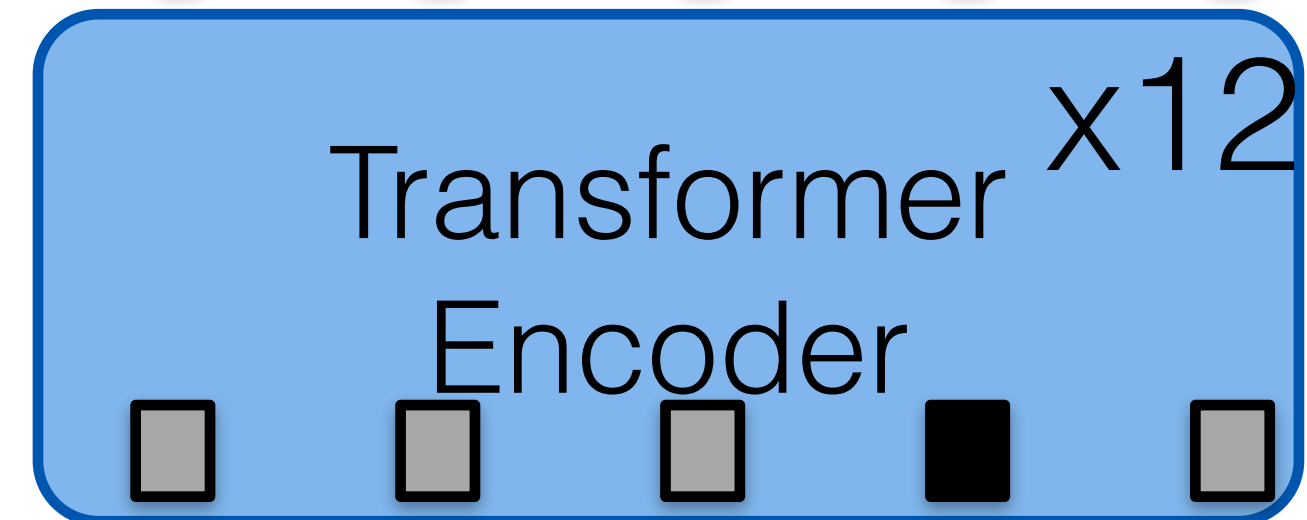
# Wav2Vec / Wav2Vec 2.0

- Architecture:
  - CNN+Transformer

- Training
  - Masked prediction of quantized vector
  - contrast true quantized latent with distractor latent embeddings

latent representation

CNN+Transformer

# Wav2Vec2

Context C

Mask during training

Quantized Rep Q

latent rep Z

Transformer Encoder ×12

CNN ×7

How many layers of Convolution?

How to design each kernel size/stride?

Raw wav X, each frame ~ 25ms, stride 20ms

# Wav2Vec2

Context C

Transformer Encoder x12

Mask during training

Quantized Rep Q

latent rep Z

Conv1D
width=2, stride=2,d=512
x2

Conv1D
width=3, stride=2,d=512
x4

Conv1D
width=10, stride=5,d=512

CNN x7

Raw wav X, each frame ~ 25ms, stride 20ms

waveform $x$ at16kHz

Wav2vec2.0: a Framework for Self-Supervised Learning of Speech Representations [Baevski et al, NeurIPS 2020]

# Wav2Vec2

Context C

Transformer Encoder x12

Mask during training

Quantized Rep Q

latent rep Z

q   q   q   q   q

CNN    x7

Raw wav X, each frame ~ 25ms, stride 20ms

frame size=399 (25ms)

sampling rate=50Hz
(sliding 320=20ms)

Conv1D
width=2, stride=2,d=512    x2

Conv1D
width=3, stride=2,d=512    x4

Conv1D
width=10, stride=5,d=512
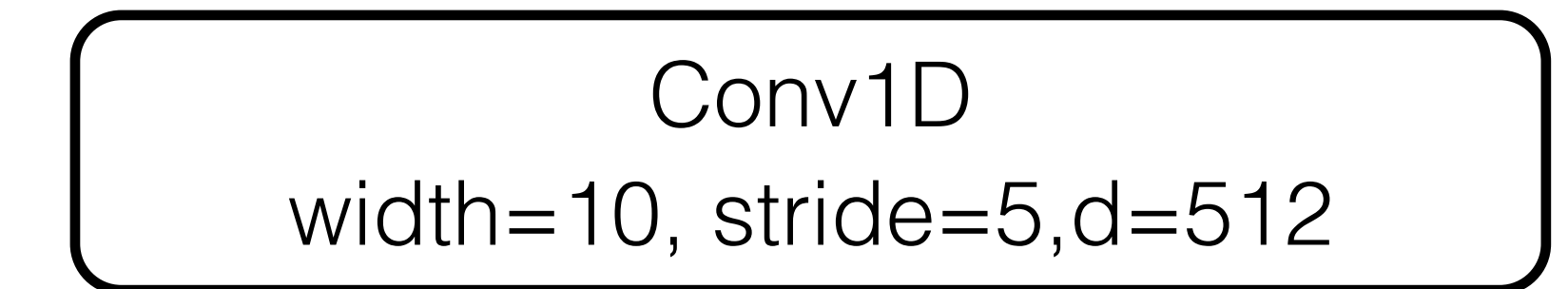
waveform $x$ at16kHz

Wav2vec2.0: a Framework for Self-Supervised Learning of Speech Representations [Baevski et al, NeurIPS 2020]    8

# Discrete Quantization with Codebook

Linear

Concat

pick max prob

G groups of probability vector of size V

| 0.1 | 0.2 | 0.3 | 0.7 |
| 0.3 | 0.3 | 0.4 | 0.1 |
| 0.6 | 0.5 | 0.3 | 0.2 |

(Gumbel) Softmax

$G \times V$

Linear

one frame vector from CNN

codebook1

codebook2

codebook3

codebookG

# How to obtain codebook — Product Quantization

Splitting a vector into equally sized chunks — subvectors,
Assigning each of these subvectors to its nearest *centroid*

# Contextual Encoder

# Wav2Vec2.0: Contrastive on quantized acoustic state



Training data: (audio only)
LibriSpeech 960 hrs
LibriVox 53k hrs

Masked context during training

**Minimize contrastive loss**

$$L = -\sum \log \frac{\exp Sim(c_t, q_t)}{\sum \exp Sim(c_t, q_-)} + \text{penalty}$$

Quantized low-level acoustic state, each frame ~ 25ms, stride 20ms

Transformer Encoder x12

CNN x7

Bring closer masked context and quantized acoustic state

12

# Training Loss



Cosine similarity · Context representation · Discrete latent speech representation

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

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 frames
  - span can overlap
  - for a 15s sample, ~49% of frames masked with an avg span of 300ms

# Model Setup

- Wav2vec2 base:
  - 12 Transformer layers, d=768, d_ffn=3072, #heads=8
  - 16 groups
  - rel pos emb cnn kernel size 128

- Wav2vec2 large:
  - 24 Transformer layers, d=1024, d_ffn=4096, #heads=16

# Training

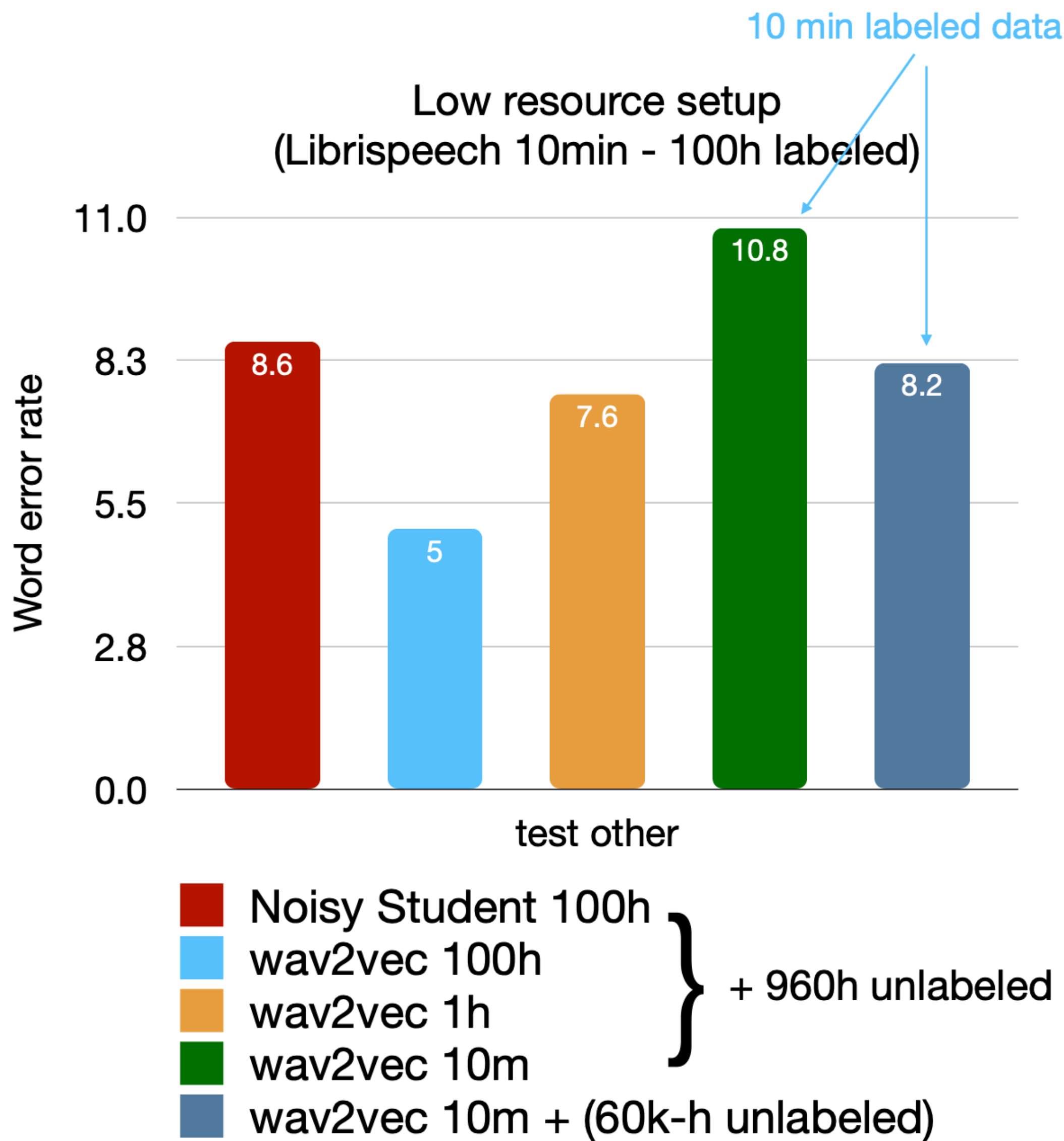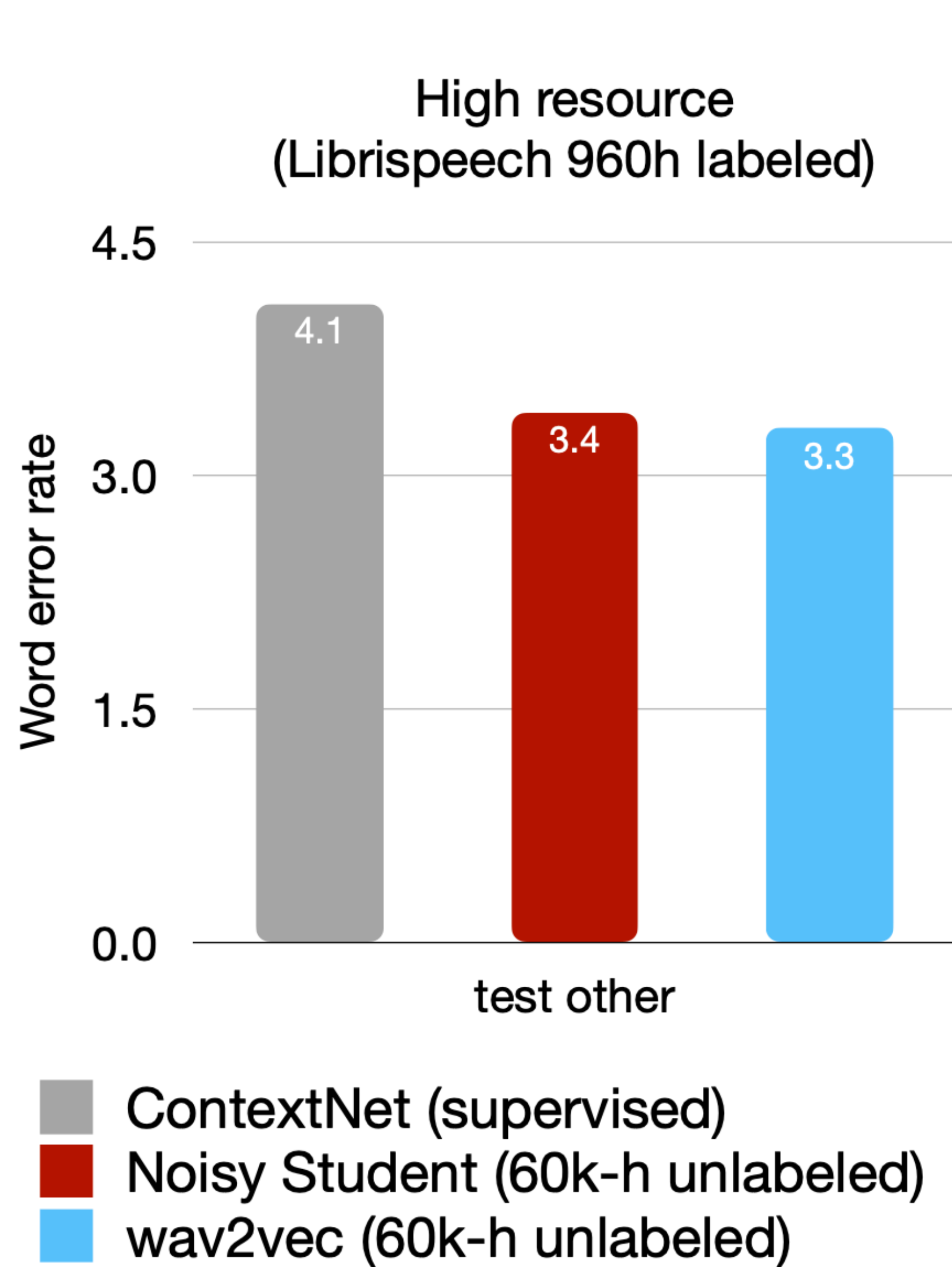- LibriSpeech: 960 hours of English speech (just audio)

- LibriVox (LV-60k): about 53k hours of audio for book reading

- Wav2Vec2 base:
  - each sample is cropped with length 250k (=15.6s)
  - total batch size: 1.6 hours on 64 V100 GPUs

- Wav2Vec2 Large:
  - each sample is cropped with length 320k (=20s)
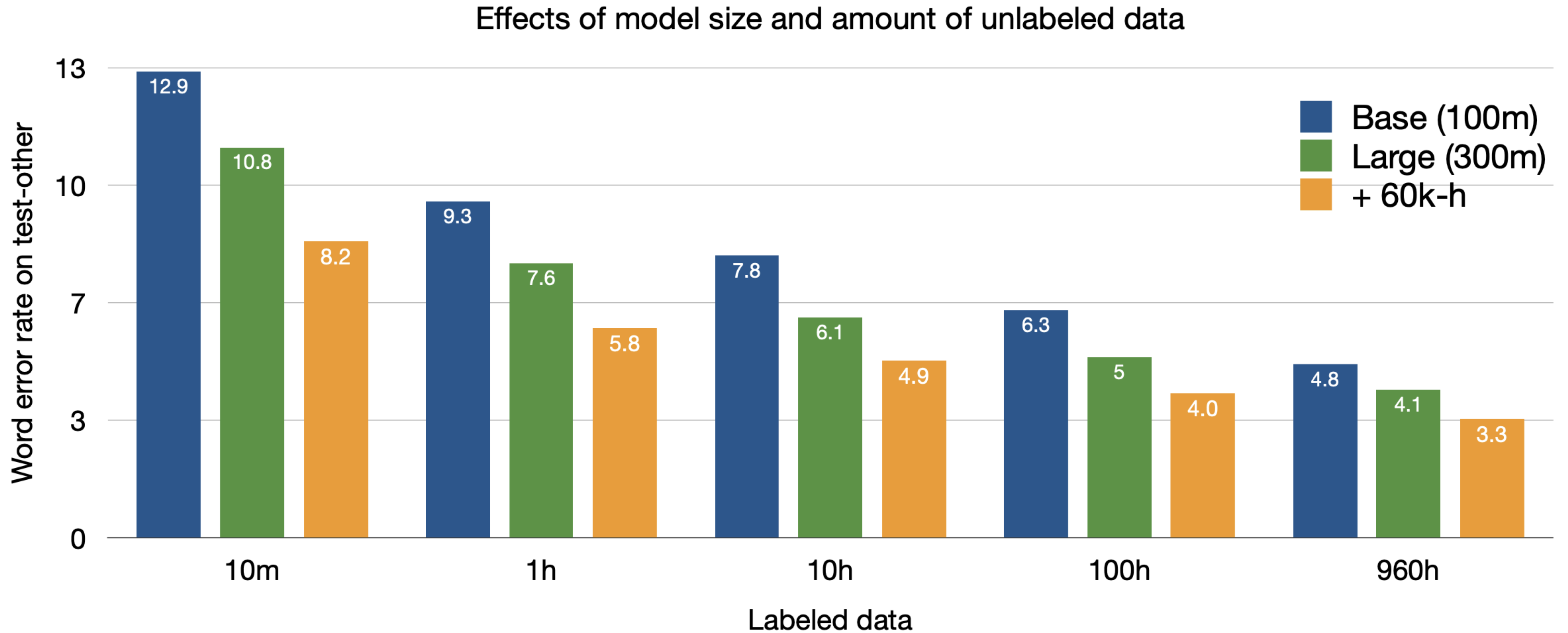  - total batch size: 2.7hours on 128 V100 GPUs.

# 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 wav to letter generation with the official 4gram LM and Transformer LM
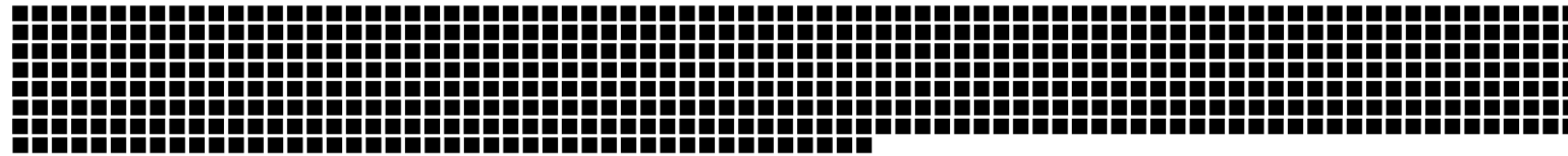
# Wav2Vec2 Results

# Effects of Model size and raw data



Effects of model size and amount of unlabeled data

Word error rate on test-other

Legend:
- Base (100m) — blue
- Large (300m) — green
- + 60k-h — orange

| Labeled data | Base (100m) | Large (300m) | + 60k-h |
|---|---|---|---|
| 10m | 12.9 | 10.8 | 8.2 |
| 1h | 9.3 | 7.6 | 5.8 |
| 10h | 7.8 | 6.1 | 4.9 |
| 100h | 6.3 | 5 | 4.0 |
| 960h | 4.8 | 4.1 | 3.3 |

# Overall ASR results



Librispeech benchmark, WER on test-other

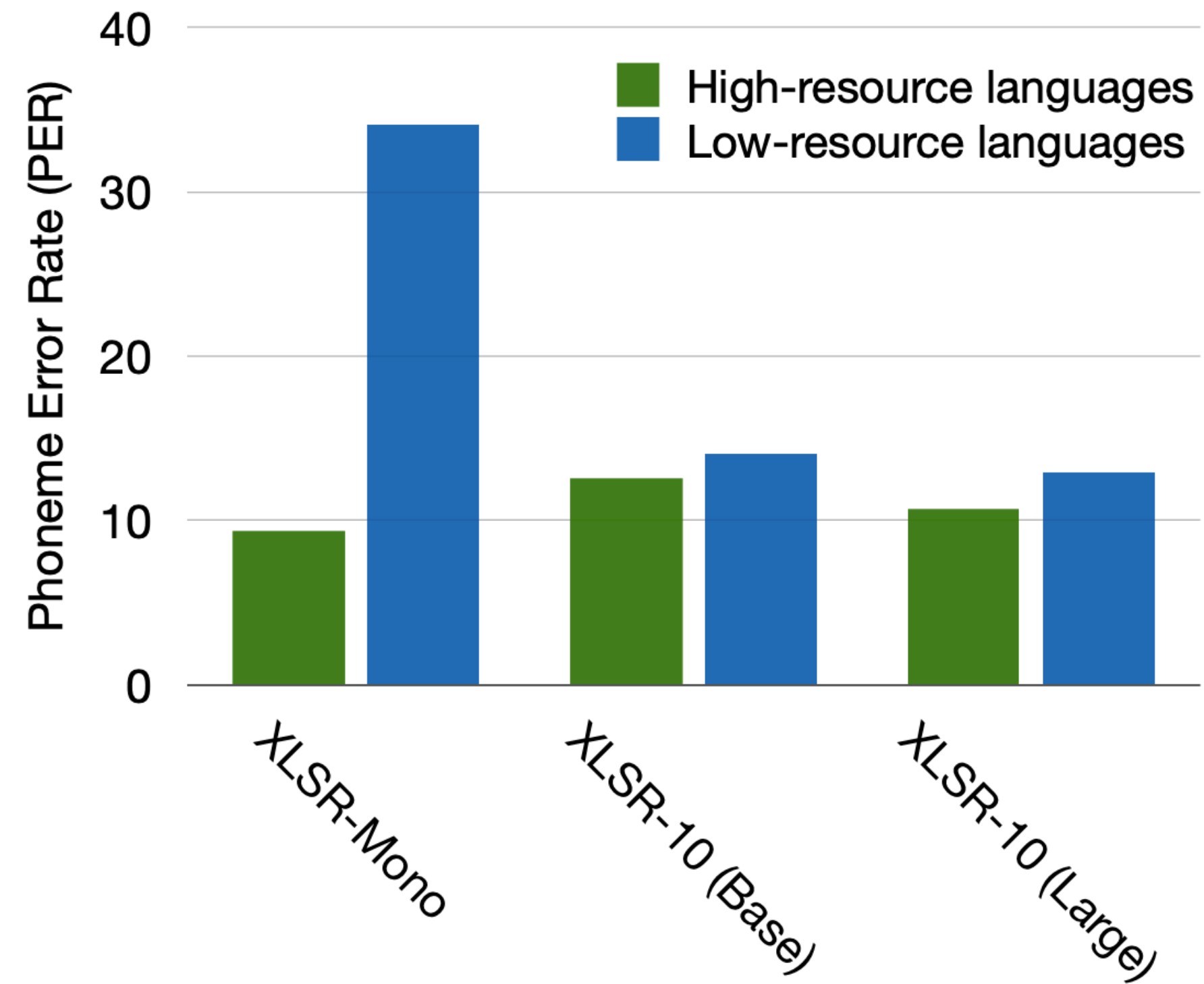Data based on Papers with Code (25 Oct 2020)
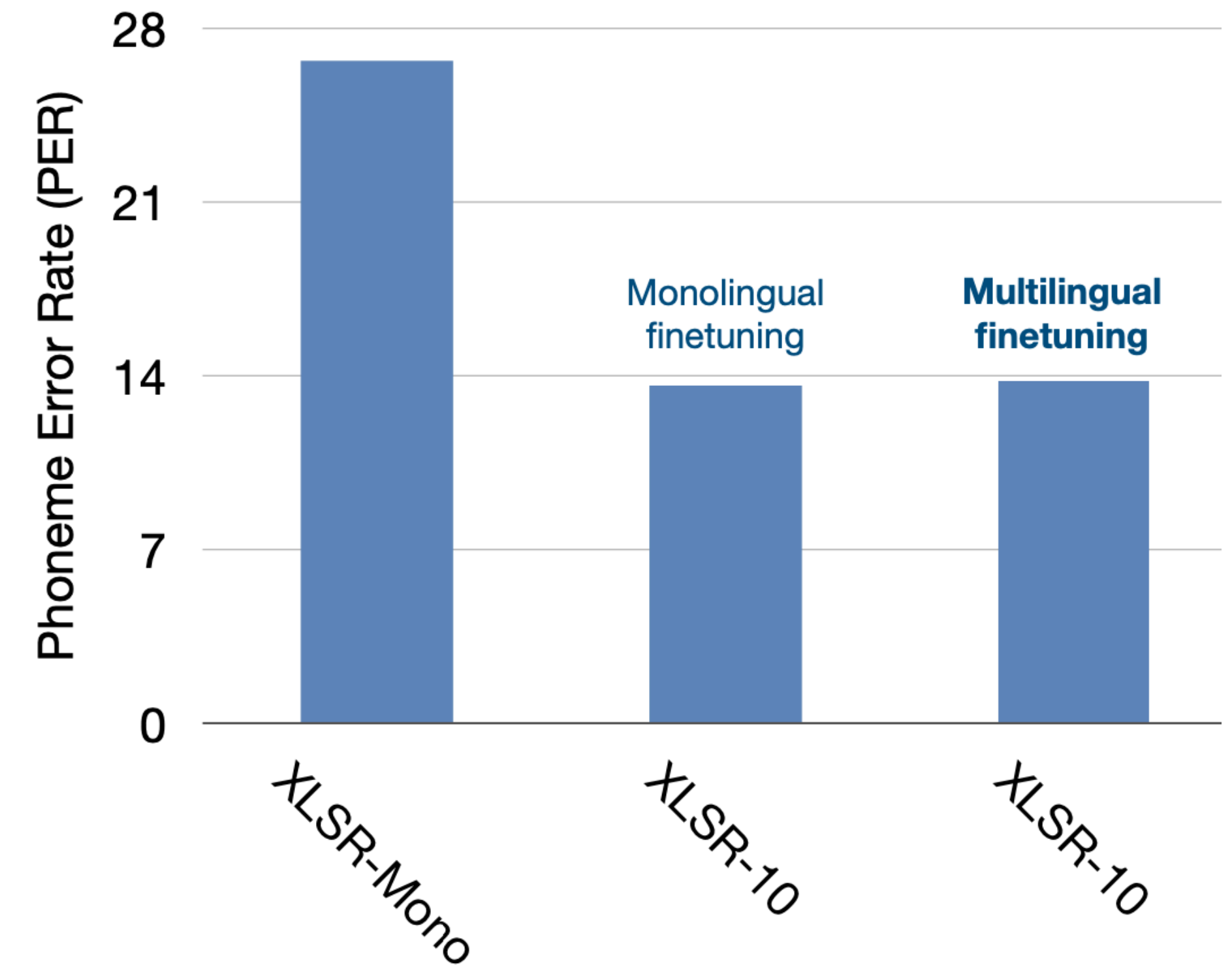
# XLSR: Multilingual Wav2Vec2

## Cross-lingual transfer

**CommonVoice results:**



## Multilingual fine-tuning

**CommonVoice results:**

# Summary

- Self-supervised pre-training with audio data only

- Wav2Vec2 Model: CNN+Transformer

- construct the frames with reasonable size (25ms) and sliding (20ms)
  - proper design of CNNs

- Masked training with contrastive loss on quantized representation

22

# Language in 10

# TTS Code in Notebook

- https://github.com/lileicc/FastSpeech2

- https://www.cs.cmu.edu/~leili/course/11737mnlp23fa/code/tts/run_tactron2.ipynb