# S3PRL introduction & recent update in JSALT

Shu-wen (Leo) Yang

Andy T. Liu, Heng-Jui Chang, Haibin Wu, Cheng Liang

# S3PRL

**S**elf-**S**upervised **S**peech **P**re-training
and **R**epresentation **L**earning

https://github.com/s3prl/s3prl/

## s3prl

s3prl

Self-Supervised Speech Pre-training and Representation Learning Toolkit.
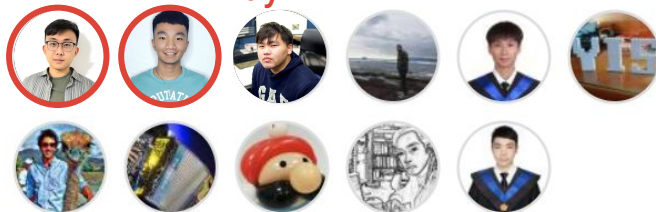
🔗 youtu.be/PkMFnS6cjAc

⭐ 1.4k stars    🍴 315 forks

https://github.com/s3prl/s3prl/

Used by 14

+ 6

Contributors 38

Creators

Leo    Andy

+ 27 contributors

Prof. Hung-yi Lee, Advisor & Sponsor

# Major functionality

No version

v0.2

v0.3

Heng-Jui Chang

Xuankai Chang

Yung-Sung Chuang

Pre-training

Pre-trained model collection

Downstream Benchmarking & SUPERB

Zili Huang

Wen-Chin Huang

Tzu-Hsien Huang

Kushal Lakhotia

Yist Lin Y.

Guan-Ting Lin

**2019**

**2020**

**2021**

Andy T. Liu

Shu-wen Yang

Shu-wen Yang

Jiatong Shi

Shu-wen Yang

Andy T. Liu

Andy T. Liu

Hsiang-Sheng Tsai

Po-Han Chi

Po-Han Chi

Wei-Cheng Tseng

# Feedback ⟵⟶ Improvement during JSALT workshop

- It is intensively used in the JSALT pre-training team for evaluating new techniques
  - Bugs reported
  - Thank all the users for reporting the error
  - Thank JSALT for providing the platform to have so many users to help the open-source

- More efficient
- Better generalization
- Visual enhanced

| How to better use SSL models | Enhance SSL models |
|---|---|
| Push SSL models to more tasks | Toolkit |

- Prosody-related Tasks
- Spoken Language Understanding

# Feedback ⟷ Improvement during JSALT workshop

- We receive lots of feedback and continuously improve it:

    - How to change the corpus for XXX task?

    - How to change the probing model for the XXX task?

    - The steps to benchmark a new SSL model is too complicated

    - Connection to the HuggingFace models

    - How to benchmark with just a subset of the corpus?

    - The latest SSL models?

S3PRL was not designed as a flexible/reusable library
but as the recipes to reproduce papers

# Feedback ←→ Improvement during JSALT workshop

- The mostly asked issue is…

## Fairseq installation issues

- At some commits, some pre-trained models work, but the others fail
- At some commits, all the pre-trained models work, but the new models can't be supported
- At some commits, you can't even `torch.load` the some checkpoints, since the checkpoints contain deprecated pickled Fairseq object

## Too Difficult…

# Feedback ⟵⟶ Improvement during JSALT workshop

- Fixing these issues → saving lots of users' time

- So the users can use S3PRL under more settings without tracing the code

s3prl

## s3prl

Self-Supervised Speech Pre-training and Representation Learning Toolkit.

🔗 youtu.be/PkMFnS6cjAc

☆ **1.4k** stars    **315** forks

# Major Updates for v0.4

- Deprecate tasks' God Classes
- Hooks for changing corpus, downstream, and upstream
- More upstream models
- Remove Fairseq dependency
- HugggingFace connection
- Audio / Sound connection

# The major cause of most of the issues

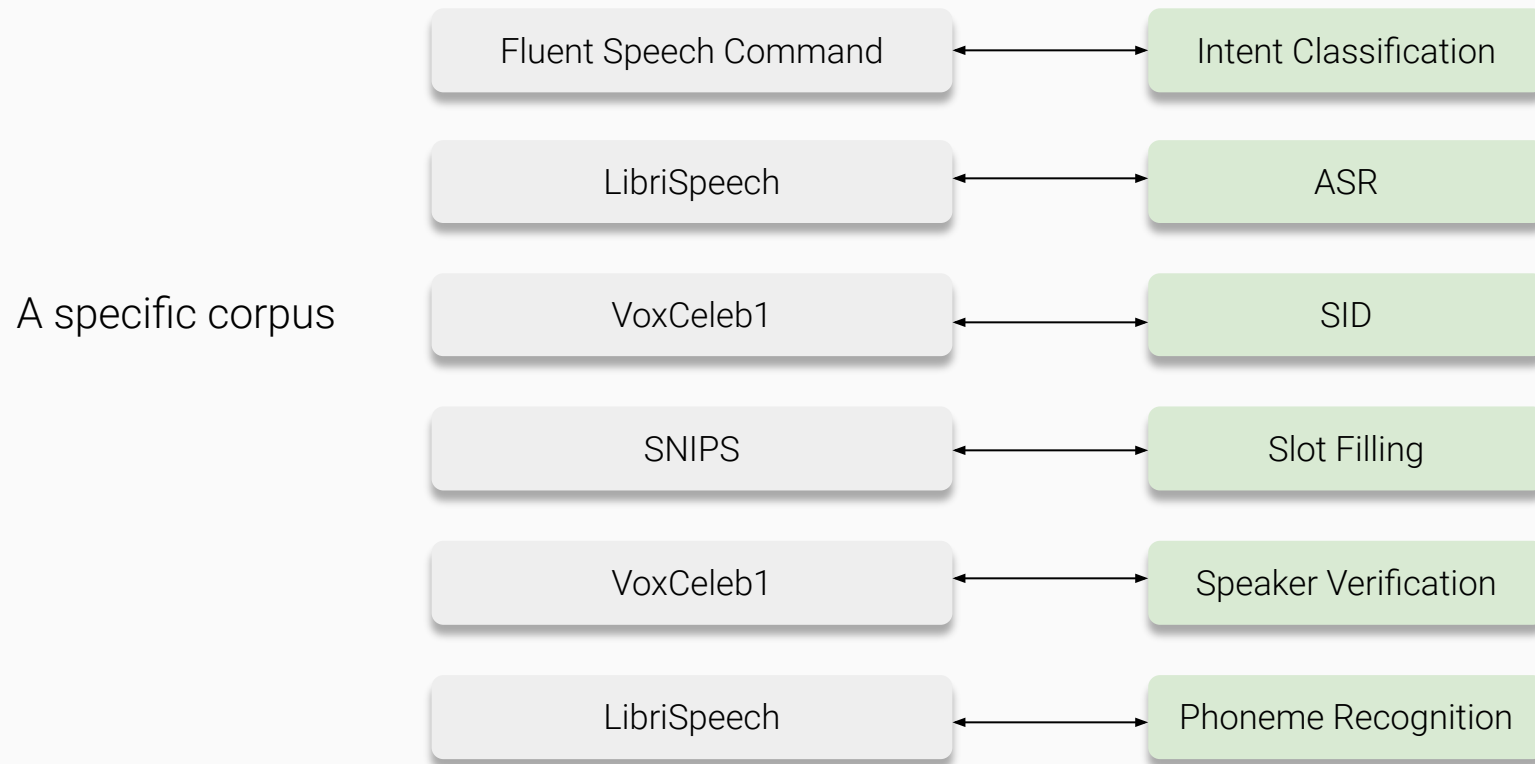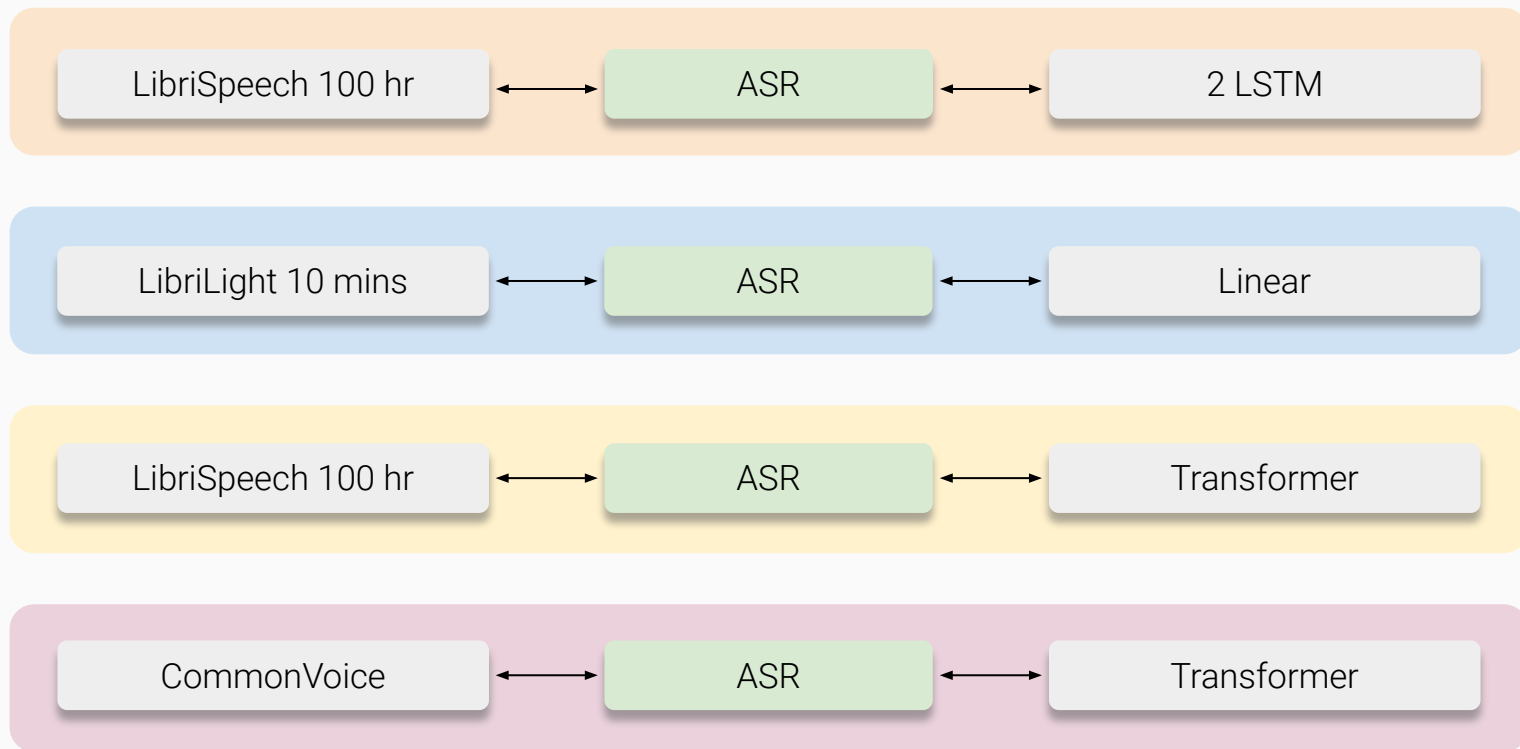- A single God Class handling all the details for a task

nn.Module

yenmeng fix ctc decode

..

a2a-vc-vctk

a2o-vc-vcc2020

asr

asr_lxt

atis

audio_snips

ctc

diarization

docs

emotion

enhancement_stft

example

fluent_commands

libri_phone

lxt_dtw

# In the God Class - Model entanglement

A specific model

| | |
|---|---|
| Mean Pooling + Linear | Intent Classification |
| 2 LSTM | ASR |
| Mean Pooling + Linear | SID |
| 2 LSTM | Slot Filling |
| Xvector (TDNN, Stats Pooling) | Speaker Verification |
| Phoneme Recognition | Phoneme Recognition |

# In the God Class - Corpus entanglement

| | | |
|---|---|---|
| **Fluent Speech Command** | ← → | Intent Classification |
| **LibriSpeech** | ← → | ASR |
| A specific corpus **VoxCeleb1** | ← → | SID |
| **SNIPS** | ← → | Slot Filling |
| **VoxCeleb1** | ← → | Speaker Verification |
| **LibriSpeech** | ← → | Phoneme Recognition |

# Every slightly change requires library code change

| LibriSpeech 100 hr | ⟷ | ASR | ⟷ | 2 LSTM |

| LibriLight 10 mins | ⟷ | ASR | ⟷ | Linear |

| LibriSpeech 100 hr | ⟷ | ASR | ⟷ | Transformer |

| CommonVoice | ⟷ | ASR | ⟷ | Transformer |

# Disentangle corpus from the task

# Hooks to customize the behavior for each task

corpus hook
Change the data

upstream hook
Change the feature model (SSL)

downstream hook
Change the task model

train/valid/test dataset hook
Change the __getitem__
(add noise, reverb… etc)

train/valid/test sampler hook
Change the batching behavior

# Default usage

- Default hooks reproduce the exact SUPERB setting

```
classmethod run(**cfg)
```

- *Necessary Config:*

```yaml
workspace: ???     # (str) The workspace shared across stages
setup:
  corpus:
    dataset_root: ???     # (str) The root path of the corpus
  upstream:
    name: ???
```

```python
20  SuperbASR.run(
21      workspace="result/pseudo_asr",
22      setup=dict(
23          corpus=newdict(
24              dataset_root="/home/leo/d/datasets/LibriSpeech",
25          ),
26          upstream=dict(
27              name="hubert",
28          ),
29      ),
30      train=dict(
31          optimizer=dict(
32              lr=1.0e-2,
33          )
34      )
35  )
```

```yaml
train:
  n_jobs: 4     # (int) The number of jobs when multiprocessing on CPU
  seed: 1337    # (int) The seed
  device: cuda:0    # (str) The device used for training
  rank: 0       # (int) The global rank when distributed training
  world_size: 1     # (int) The total number of processes when distributed training
  optimizer:
    CLS: torch.optim.adam.Adam     # (str) The class used to create the optimizer. The below
    lr: 0.0001
```

# Hooks to customize data

```
setup:
  corpus:
    CLS: librispeech_for_speech2text    # (str)
  # The corpus class. You can add the **kwargs right below this CLS key
    dataset_root: ???    # (str) The root path of the corpus
```

```python
1    from pathlib import Path
2    from s3prl import newdict
3    from s3prl.problem import SuperbASR
4    from s3prl.util.pseudo_data import pseudo_audio
5
6    N_SAMPLES, N_TRAIN, N_VALID, N_TEST = 40, 20, 10, 10
7
8    def prepare_pseudo_data(wav_paths):
9        train_paths = wav_paths[:N_TRAIN]
10       valid_paths = wav_paths[N_TRAIN : N_TRAIN + N_VALID]
11       test_paths = wav_paths[N_TRAIN + N_VALID :]
12
13       def path_to_datapoint(path):
14           return {
15               "wav_path": path,
16               "transcription": "Hello World",
17           }
18
19       train_data = {Path(path).stem: path_to_datapoint(path) for path in train_paths}
20       valid_data = {Path(path).stem: path_to_datapoint(path) for path in valid_paths}
21       test_data = {Path(path).stem: path_to_datapoint(path) for path in test_paths}
22
23       return {
24           "train_data": train_data,
25           "valid_data": valid_data,
26           "test_data": test_data,
27       }
```

```python
with pseudo_audio(secs=range(1, N_SAMPLES + 1), sample_rate=16000) as (
    wav_paths,
    num_samples,
):

    SuperbASR.run(
        workspace="result/pseudo_asr",
        setup=dict(
            corpus=newdict(
                CLS=prepare_pseudo_data,  # hook for changing data
                wav_paths=wav_paths,  # arguments to the hook
            ),
            upstream=dict(name="hubert"),
        ),
    )
```

**Note:** It is also easy to load Kaldi based data directory by using a directory parser hook

# Hooks to customize downstream

```python
import torch.nn as nn
from s3prl import newdict
from s3prl.problem import SuperbASR

class CustomizedModel(nn.Module):
    def __init__(self, input_size, output_size, hidden_size: int) -> None:
        super().__init__()
        self.model = nn.Sequential(
            nn.Linear(input_size, hidden_size),
            nn.Linear(hidden_size, output_size),
        )

    def forward(self, x, x_len):
        x = self.model(x)
        return x, x_len

SuperbASR.run(
    workspace="result/pseudo_asr",
    setup=dict(
        corpus=dict(
            dataset_root="/home/leo/d/datasets/LibriSpeech",
        ),
        upstream=dict(
            name="hubert",
        ),
        downstream=newdict(
            CLS=CustomizedModel,
            hidden_size=256,
        ),
    ),
)
```

# Hooks to customize upstream

```python
import torch
import torch.nn as nn
from s3prl import newdict
from s3prl.problem import SuperbASR

class CustomizedUpstream(nn.Module):
    def __init__(self, ckpt_path: str) -> None:
        super().__init__()
        ckpt = torch.load(ckpt_path, map_location="cpu")
        hidden_size = ckpt["config"]["hidden_size"]
        self.model = nn.Sequential(
            nn.Linear(1, hidden_size),
            nn.Linear(hidden_size, hidden_size),
        )
        self.model.load_state_dict(ckpt["model_weights"])

    def forward(self, x, x_len):
        x = self.model(x)
        return x, x_len

SuperbASR.run(
    workspace="result/pseudo_asr",
    setup=dict(
        corpus=dict(
            dataset_root="/home/leo/d/datasets/LibriSpeech",
        ),
        upstream=newdict(
            CLS=CustomizedUpstream,
            ckpt_path="./ckpts/ssl_val_best.ckpt",
        ),
    ),
)
```

# More upstream models

## In SUPERB

| | | | | | | |
|---|---|---|---|---|---|---|
| Mockingjay | TERA | HuBERT | APC | VQ-APC | NPC | DeCoAR |
| DeCoAR 2.0 | Modified CPC | wav2vec | vq-wav2vec | wav2vec 2.0 | PASE+ | |

## New models

| | | | |
|---|---|---|---|
| discreteBERT | HuBERT-MGR | LightHuBERT | FitHuBERT |
| Unispeech-SAT | WavLM | data2vec | AudioAlbert | DistilHuBERT |

# To give a more stable support for users and ESPNet

- **Remove all the fairseq dependencies**
  - All the upstream can be used without installing fairseq
- **Add unit-tests for the forward and backward for all upstreams**
  - Test the representation and gradient's numerical values
  - Guarantee the same representation across S3PRL versions

# Reproduced results

- A complete re-build to
  - Get away all the old dirty design at once
  - Ensure we always have the exact old codebase for SUPERB for reproducibllity
  - Results on Hubert Base

| Task | PR | IC | SID | KS | ER | ASR | QBE | SD | SF | SV |
|------|------|--------|-------|-------|-------|------|------|------|-------|------|
| Metric | PER | ACC | ACC | ACC | ACC | WER | MTWV | DER | CER | EER |
| Old | 5.41 | 98.34 | 81.19 | 96.3 | 64.92 | 6.11 | 7.37 | 5.88 | 25.2 | 5.11 |
| New | 5.483 | 98.207 | 80.69 | 96.62 | 64.76 | 6.14 | 7.37 | 5.8 | 24.22 | 5.15 |
| Relative | -1% | -0.1% | -0.6% | 0.3% | -0.2% | -0.4% | 0% | 1.3% | 3% | -0.7% |

# For the last week in JSALT

- Connection to HuggingFace
- Connection to the HEAR Benchmark

# Connection to HuggingFace

- There are fewer model architectures available in Huggingface

  - wav2vec 2.0, HuBERT, data2vec, UniSpeech, WavLM

- But a lot more pre-trained checkpoints available, e.g. LeBenchmark

# Connection to Audio SSL

# Connection to HEAR Benchmark

- Transferbility of the speech SSL SOTA to the more sounds

- Exisiting HEAR codebase: <span style="color:red">official, faster</span>

  - Dump a single layer frozen representation

  - Train the downstream model

- S3PRL, on-the-fly feature extraction: <span style="color:red">slower, more flexible</span>

  - Enable examination for lots of speech SSL models

  - Enable weighted-sum over all layers

  - Enable finetuning SSL models on audio tasks

# Connection to HEAR Benchmark

**HEAR Benchmark**   Tasks   Code   Paper   Leaderboard   API   Submit   PMLR

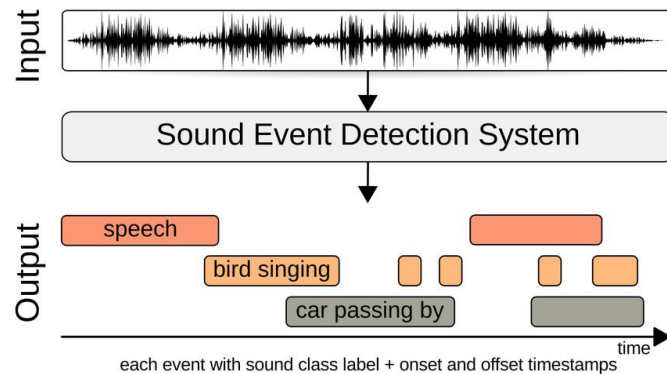| ▲ | Task Name | Embed Type | Predictor Type | Split Method | Duration (sec) | # clips | Evaluation Metric | Novel |
|---|---|---|---|---|---|---|---|---|
| ⊕ | DCASE 2016 Task 2 | T | L | TVT | 120.0 | 72 | Onset FMS | ✓ |
| ⊕ | NSynth Pitch 5hr | S | C | TVT | 4.0 | 5000 | Pitch Acc. | ✓ |
| ⊕ | NSynth Pitch 50hr | S | C | TVT | 4.0 | 49060 | Pitch Acc. | ✓ |
| ⊕ | Speech Commands 5hr | S | C | TVT | 1.0 | 22890 | Accuracy | ✓ |
| ⊕ | Speech Commands Full | S | C | TVT | 1.0 | 100503 | Accuracy | |
| ⊕ | Beehive States | S | C | TVT | 600.0 | 576 | AUCROC | |
| ⊕ | Beijing Opera Percussion | S | C | 5-fold | 4.77 | 236 | Accuracy | ✓ |
| ⊕ | CREMA-D | S | C | 5-fold | 5.0 | 7438 | Accuracy | |
| ⊕ | ESC-50 | S | C | 5-fold | 5.0 | 2000 | Accuracy | |
| ⊕ | FSD50K | S | L | TVT | 0.3-30.0 | 51185 | mAP | |
| ⊕ | Gunshot Triangulation | S | C | 7-fold | 1.5 | 88 | Accuracy | ✓ |
| ⊕ | GTZAN Genre | S | C | 10-fold | 30.0 | 1000 | Accuracy | |

# Connection to HEAR Benchmark

- **The best demonstrastion on the benefit of codebase refactoring**

- **11** new audio tasks in HEAR are immediately supported

- Task design:  k-fold, accuracy

  - Reuse the template of IEMOCAP emotion classification in SUPERB

- Change the hook configuration:

  - Corpus hook

  - Downstream model hook

# Connection to HEAR Benchmark

- 2 new tasks required to be supported in HEAR Benchmark
  - multilabel classifcation (WIP)
  - sound event detection (done)

| Rank | Model | Onset FMS |
|------|-------|-----------|
| 1 | **PaSST 2lvl+mel** | 0.9254 |
| 2 | **PaSST 2lvl** | 0.9132 |
| 3 | **wav2vec2 WS (S3PRL)** | 0.8641 |
| 12 | **wav2vec2 baseline** | 0.6630 |
| | **wav2vec2 baseline (S3PRL)** | 0.6624 |



each event with sound class label + onset and offset timestamps

# Release v0.4.0 in JSALT

| no version | v0.2 | v0.3 | | v0.4 |
|---|---|---|---|---|
| Pre-training | Pre-trained model collection | SUPERB | Heng-Jui Chang<br>Xuankai Chang<br>Yung-Sung Chuang<br>Zili Huang<br>Wen-Chin Huang<br>Tzu-Hsien Huang<br>Kushal Lakhotia | New package & Connect to Audio |
| **2019** | **2020** | **2021** | | **2022** |
| Andy T. Liu<br>Leo Yang<br>Po-Han Chi | Leo Yang<br>Andy T. Liu | Leo Yang<br>Andy T. Liu<br>Po-Han Chi | Yist Lin Y.<br>Guan-Ting Lin<br>Jiatong Shi<br>Hsiang-Sheng Tsai<br>Wei-Cheng Tseng | Leo Yang<br>Andy T. Liu<br>Heng-Jui Chang<br>Haibin Wu<br>Liang Cheng |

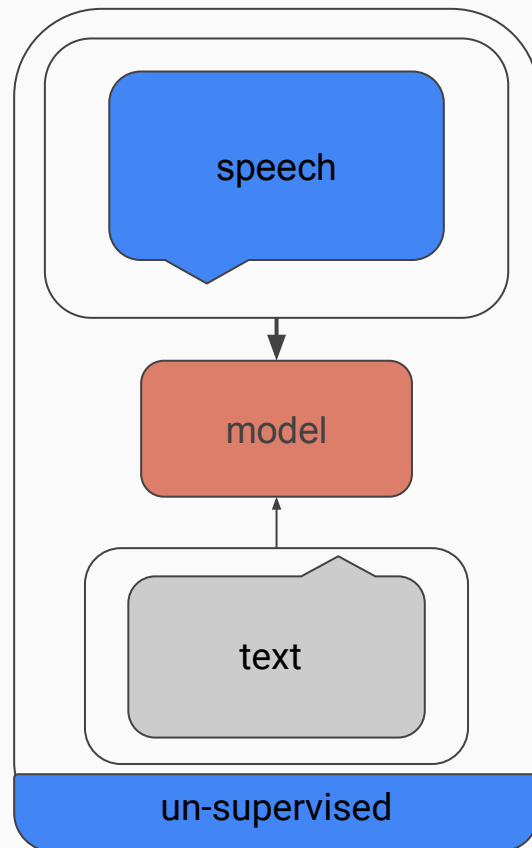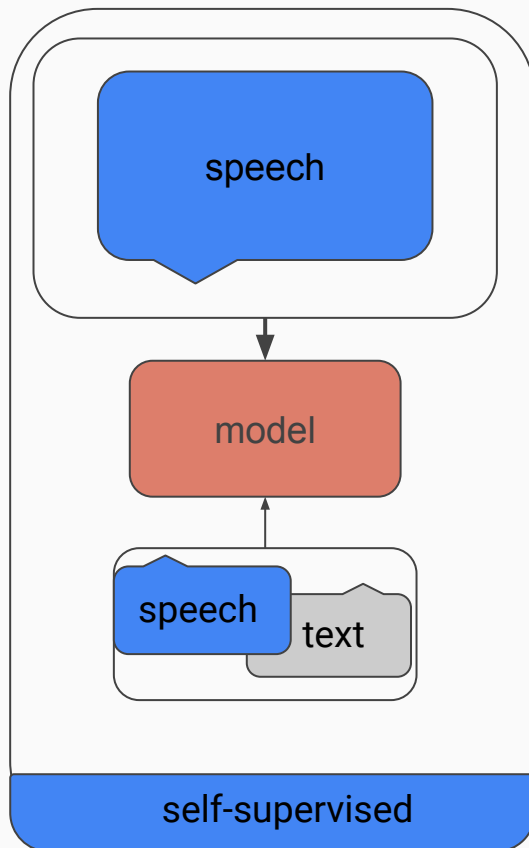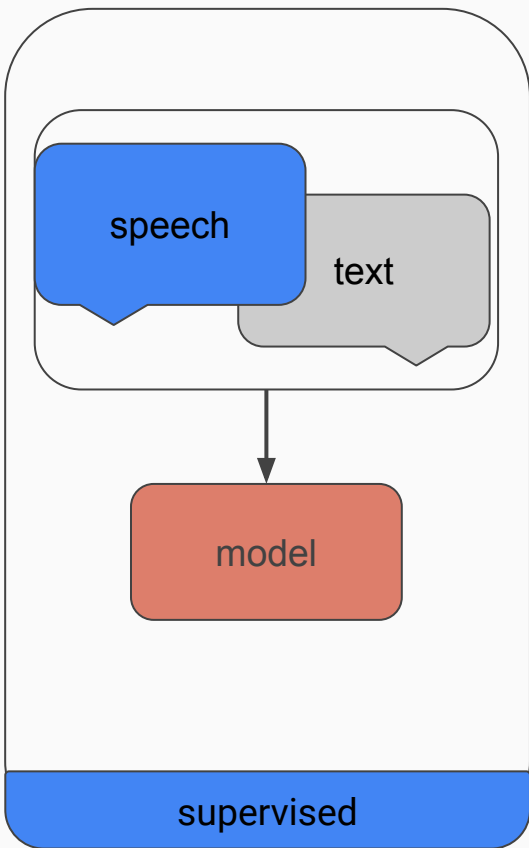# Unsupervised Automatic Speech Recognition

Dongji Gao
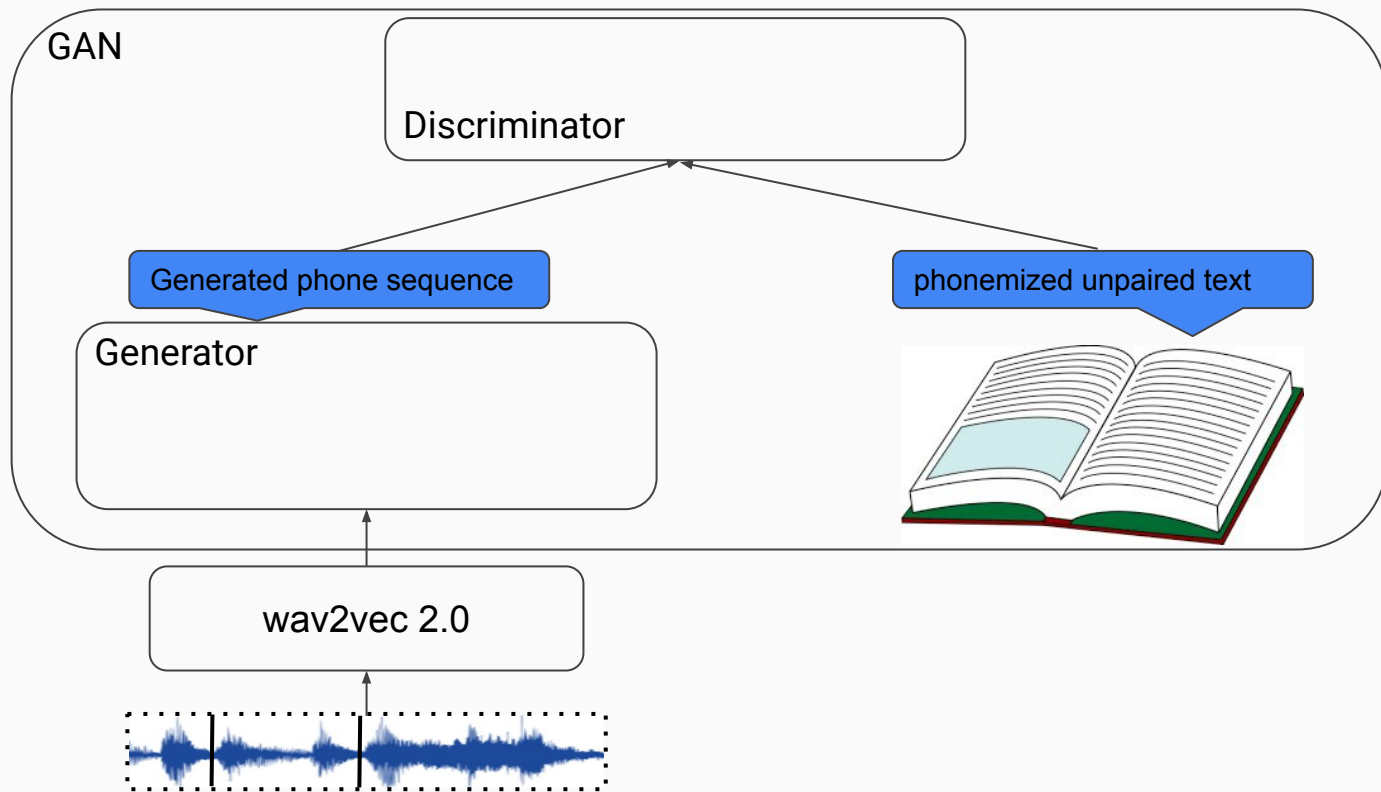
# Supervised -> self-supervised -> unsupervised

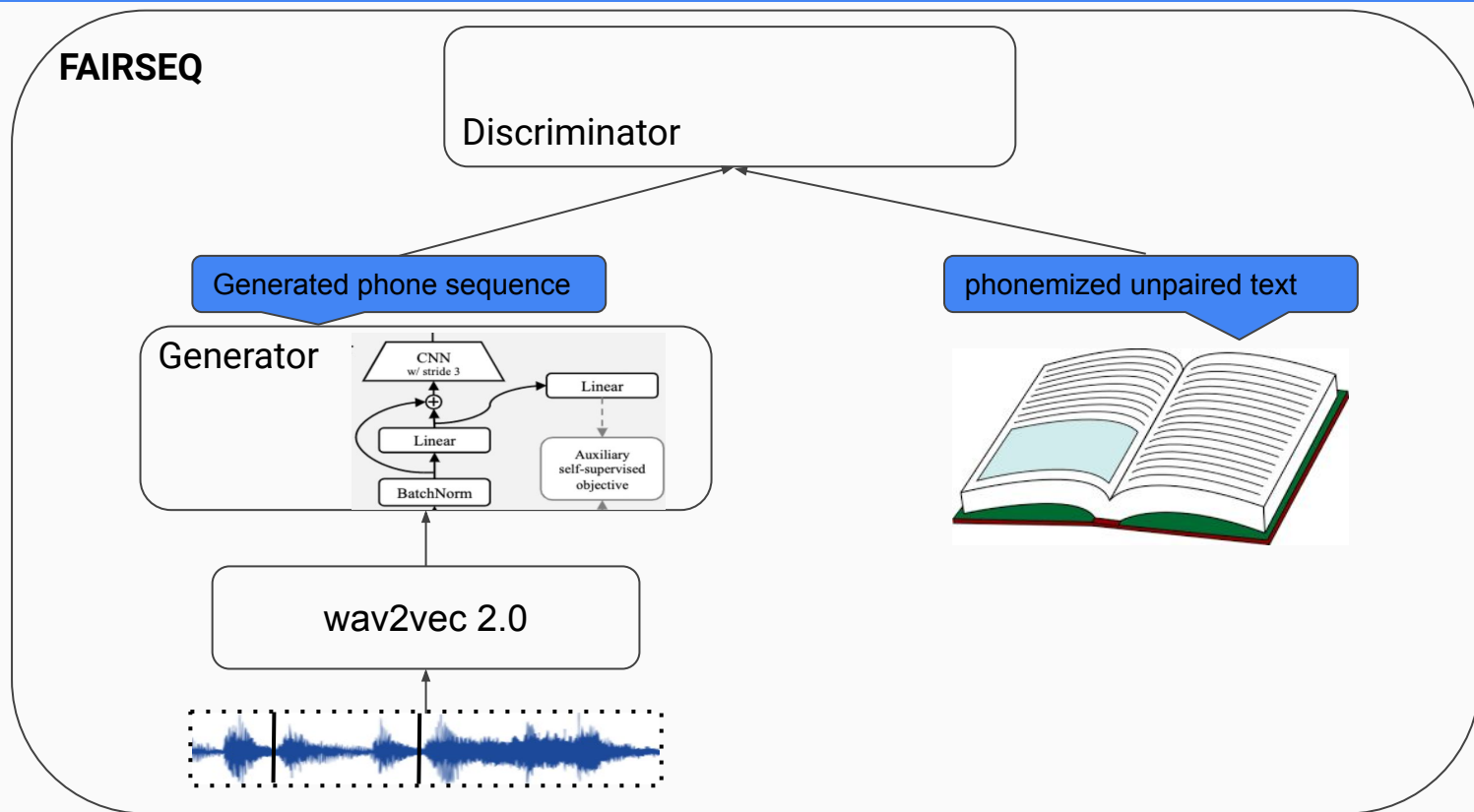# Supervised -> self-supervised -> unsupervised

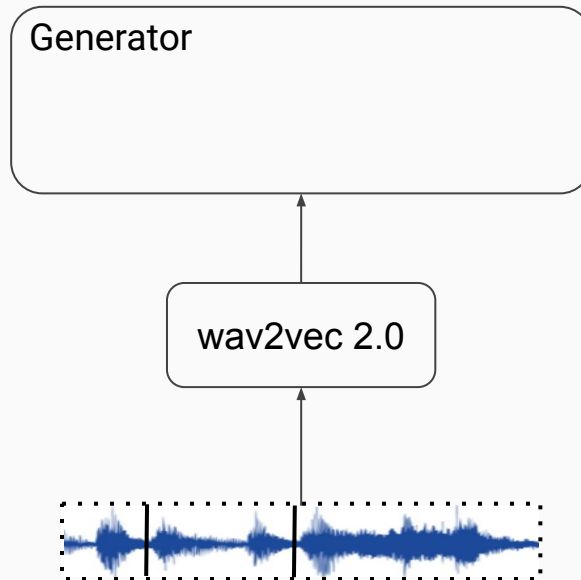# Supervised -> self-supervised -> unsupervised

# wav2vec−u2

# wav2vec−u2



FAIRSEQ

Discriminator

Generated phone sequence

phonemized unpaired text

Generator

CNN w/ stride 3

Linear

⊕

Linear

Auxiliary self-supervised objective

BatchNorm

wav2vec 2.0

Liu, Alexander H., Wei-Ning Hsu, Michael Auli, and Alexei Baevski. "Towards End-to-end Unsupervised Speech Recognition." (2022).

# ESPNET



- Front end

Generator

wav2vec 2.0

# ESPNET



- Front end

Discriminator

Generator

Hubert  wav2vec 2.0  WavLM

# ESPNET

- Front end



Generator

Hubert  |  wav2vec 2.0  |  WavLM

# ESPNET

- Front end

# ESPNET

- Front end



- Faster data preprocessing
  - Parallel
    - VAD
    - Remove silence
    - MFCC clustering
  - On-the-fly feature extraction
    - Trainable weighted sum of features from different layer

# ESPNET

- Front end



- Faster data preprocessing
  - Parallel
    - VAD
    - Remove silence
    - MFCC clustering
  - On-the-fly feature extraction
    - Trainable weighted sum of features from different layer
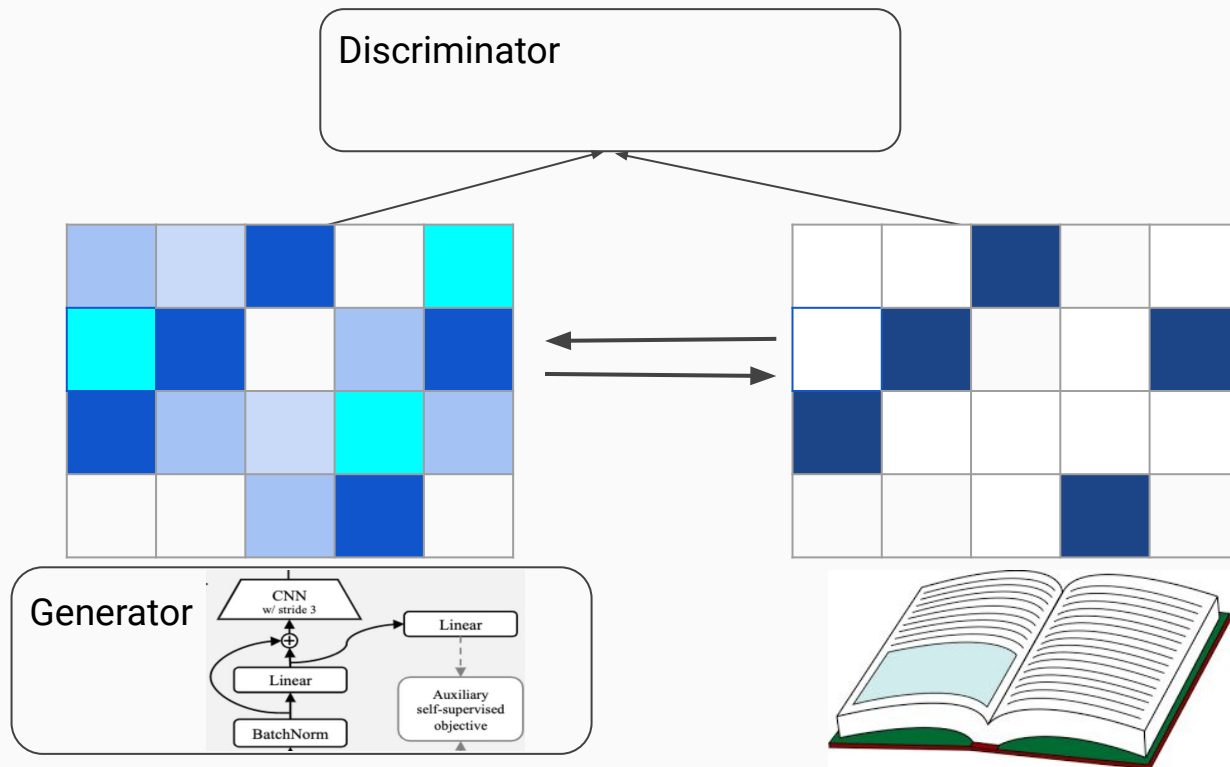- Training
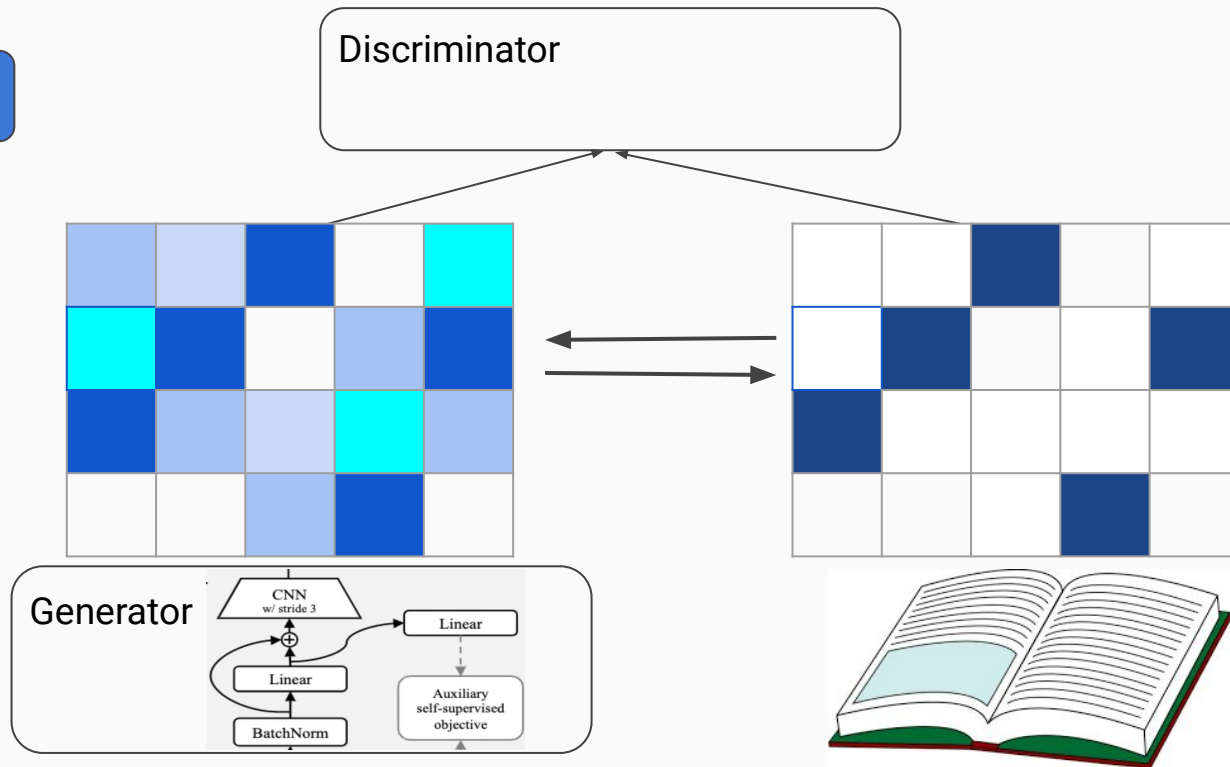  - Reproducibility
  - Efficiency
  - Performance

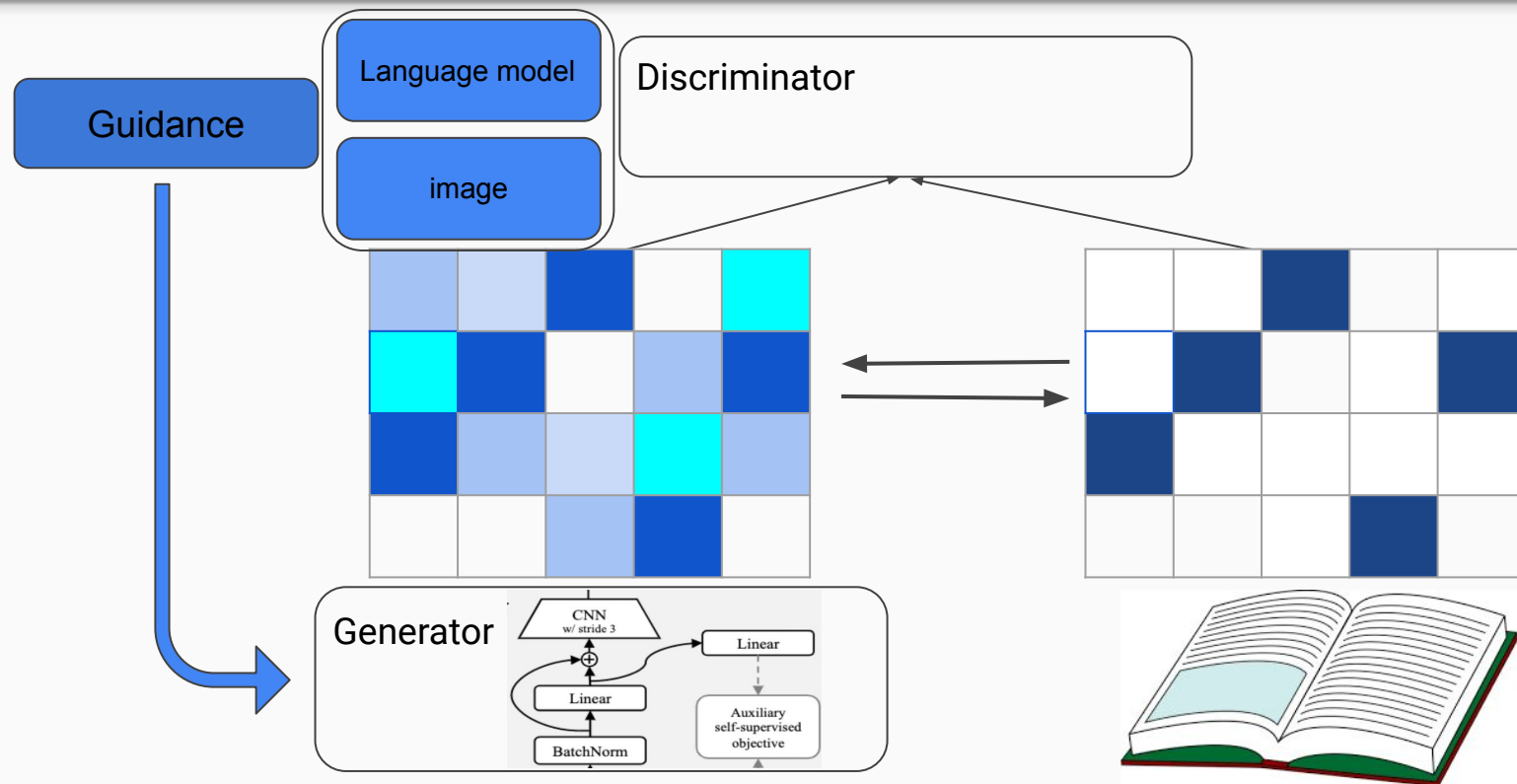Text: "C B A D B", phoneme set {A, B, C, D}

Text: "C B A D B", phoneme set {A, B, C, D}

Text: "C B A D B", phoneme set {A, B, C, D}

# ESPNET



- Front end



- Faster data preprocessing
  - Parallel
    - VAD
    - Remove silence
    - MFCC clustering
  - On-the-fly feature extraction
- Training

# Thanks!

Dongji Gao
Jiatong Shi
Ann Lee
Paola Garcia
Shinji Watanabe
Hung-yi Lee

# Results on Librispeech

| Pseudo Label | # of class | Avg. PER |
|---|---|---|
| None | - | $15.9 \pm 1.1$ |
| wav2vec2.0 VQ indices[2] | $320 \times 2$ | $16.6 \pm 2.2$ |
| k-means clustering wav2vec2.0 features | 32 | $16.4 \pm 1.4$ |
| | 64 | $15.5 \pm 1.8$ |
| | 128 | $15.9 \pm 0.9$ |
| k-means clustering MFCC audio features | 50 | $15.2 \pm 0.9$ |
| | 64 | $\mathbf{13.6 \pm 0.9}$ |
| | 100 | $14.8 \pm 1.3$ |
| | 128 | $16.8 \pm 1.7$ |

Text: "C B A D B", phoneme set {A, B, C, D}