Sora Nagano

The University of Tokyo

2025-07-12

# 1. □□□□□□□□

□ □□

□□□□□□□□□□□□□□□□□□□□□□□□□□□□

- □□□□□□□□Chomsky & Halle 1968□
- □□□□□□Prince & Smolensky 1993□
- □□□□□□□□□□□□□□□□□□□□

□□□□□□□□□□□□[sakɯᵝa] → /sakura/

- □□□/s/, /a/, /k/, /u/, /r/, /a/□□□□□□□□□
- □□□[+consonantal], [-voice], [+coronal]……□□□□□□□□□

□ □□□□□

□□□□□□□□□□□□□□□□□□□□□□□□□□□ (Staples & Graves, 2020)

- word2vec, Skip-gram□Mikolov et al. 2013□
- □□□□□□□king - man + woman ≈ queen (Silfverberg et al., 2018)
- □□□□□□□□□□□□□□ (Kolachina & Magyar, 2019)

□□□□□□□□□□□□□□

> □ □□
>
> **GAN** □□□□□□□ (Begǔs, 2020)□
> - □□□□□□□□□□□□□□□□□□□□□
> - VOT□□□□□□□□□□□□□□□□□□
> - □□□□□□□□□□□□□□□□□□□□□□□□ (Chen & Elsner, 2023)

□□□□□□□□□□□□□□

> **□ □□□□□**
>
> **Vector Quantization (VQ)**□□□ → □□□□□□□□□ (Higy et al., 2021)
> - k-means □□□□□□□ → □□□□□□□□
> - Gumbel-Softmax□□□□□□□□□
> - □□□□□□□□□□□□□□□□128, 256, 512...□□□□□□

> **□ □□**
>
> VQ □□□□□
> □□□□□□[0.3, 0.8, −0.2] → □□□□□□□□ID:47□
> □□□□□□□□□□□□□□□□□□□□□□□□□

**SSL □□□□□□□□□**

> □ □□
>
> □□□□□□□□□□□□ (Venkateswaran et al., 2025)□
> - □□□□□□□□□□□□□□□F0, formant□
> - □□□□□□□□□□□□□□□□□
> - □□□□□□□□□□□□□□□□□

□□□□□□□

> □ □□□□□
>
> - □□□□□□□□/p/-/pʰ/□□□ (Medin et al., 2024)
> - □□□□□□□□□□□□□□□ (Pasad et al., 2024)
> - □□□□□□□□□□□□□□□□□□□ (Pouw et al., 2024)
> - □□□□□□□□□□□□□□□□□ (Gosztolya et al., 2024)

□□□□□□□□□□□□□

> **□ □□□□□**
>
> □□×□□□□□□□□□□□ (Panchendrarajan & Zubiaga, 2024)□
>
> - □□□□□□□□+□□□□□□□□□
> - □□□□□□□□□□□□□□□□□□□

> **□ □□**
>
> □□□□□□□□□□□
>
> □□□□□□→□□□□□□□□□
>
> □□□→□□□□□□□□□
>
> □□□□□□□NN□→□□□□□□□→□□□□

# 1.5 □□□□□□□

□□□□□□□

□□□□□□□□□□□□□□□□□□□□□□

- □□□□□□□□□□ $X$ □□□ $Y$ □□□□□
- □□□□□□□□□□□□□□□□□□
- □□□□□□□□□□□□□□□□□□

□□□□□□□

> **□ □□**
>
> □□□□□□□□□□□□□□□□
>
> - □□□□□□□□VQ□□□□□□□□□□
> - □□□□□□□□□□□□□□□□□□□□
> - □□□□□□□□□□□□□□□□□□□□□□□□□□□□

## 2. □□□□□

# 2.1 □□□□□

□□□□□□□□□□□□□□□□□□□□□□□□

> **□ □□**
>
> □□□□□□□□□□□□□□□□□□□□ (Cho et al., 2025)□vs □□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□

- □□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□
- □□□□□□□□□□□□□□□□□□□□□□□□□□□□□□ (Staples & Graves, 2020)

# 2.2 NLP □□□□

> **□ □□**
>
> NLP = Natural Language Processing□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□
> □□Google □□□□Siri□ChatGPT

## NLP □□□□□□

- **1950 □**□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□
- **1990 □**□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□
- **2010 □**□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□

> **□ □□□□□**
>
> □□□□□ =□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□

# 2.3 □□□□ NLP □□□

## □□□□□□□□□

- □□□□□□□□□□□□□□□□□□□□□□□
- □□□□□□□□□□□□□□□□
- □□□□□□□□□□□□□□

## NLP □□□□□□

- □□□□□□□□□□□□□□□□□□□□□
- □□□□□□□□□□□□□□□□
- □□□□□□□□□□□□□□□□□

> **□ □□**
>
> □□□□□/p/□/b/□[labial]□[± voice]□□□□□□□□→ NLP □□□□□□□□□□□□□□□□$[0.2, -0.8, 1.3, \ldots]$□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□

# 3. SSL□□□□□□□□□□□□

□ □□

□□□□□□□□□□□**Self-Supervised Learning, SSL**□ =□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□
□□

□□□□□□□/a/□□□□□/k/□□□□□□□□□□□□□
⇒ SSL□□□□□□□□□□□□□□□□□□□□□□□□□ (Mohamed et al., 2022)□

# 3.1 SSL□□□□□□□□□□□□□

## □□□□□□□

- □□□□□□□□□□□□□□□□□□□□
- □□□□□□□□□□□□□□□□□□□□□□
- □□□□□□□□□□□

## SSL □□□□

- □□□□□□□□□□□□□□□□□□□□□
- □□□□□□□□□□□□□□□□□ (Choi et al., 2024)
- □□□□□□□□□□□□□□□□

**□ □□□□□**

**wav2vec 2.0**□Meta/Facebook □□□= □□□□□□□□□□□□ SSL □□□ (Baevski et al., 2022)□

- □□□□□□□LibriSpeech□960 □□□□□□□□□□□□
- □□□□□□□□□Transformer□□□□□□□□□□□□□□□□□□
- □□□□□□□□□□□□□□□□□□□□□□□

3. SSL□□□□□□□□□□

# 3.2 wav2vec 2.0□□□□□ SSL □□□

□□□□□□

1. □□□□□□□□□□□□□□□□□□□□□□□□
2. □□□□□□□□□□□□□□□□□□□□□□□
3. □□□□□□□□□□□□□□
4. □□□□□□□□□□□□□□□□□□□□□□□

> □ □□
>
> □□□□□□ __ □□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□→□□□□□□□□□□□

> □ □□
>
> □□□□□□□ =□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□
> □□□ (Astrach & Pinter, 2025; Venkateswaran et al., 2025)□

# 3.3 □□□□□□□□□□□□□□□□□

## □□□□□□□□□□

- □□□□□□□□□□□□□□□□□□□□□□□□
- □□□□□□□□□□□□□□□□□□□□□□□□
- □□□=□□□□□□□□□□□□□

# 4. VQ□□□ ⟺ □□□□□□

> □ □□
>
> **□□□□□□□□□Vector Quantization, VQ□□**
> = □□□□□□□□ ⟶ □□□□□□□□□□□□□□□□□□□□□□
>
> □□□□□[0.73, −0.45, 1.23] ⟶ □□□□□□□ID:15□
> ⟹ NLP □□□□□□□□□□□□□□□□□□□□□□□□□□ (Higy et al., 2021)□

**□□□□□□**

- □□□□□□□□□□□□□□□□□/p/, /t/, /k/□□□
- NLP□□□□□□□□□□□□
- VQ□□□□□□□□□

# 4.2 K-means □□□□ VQ □□ □□□□

1. wav2vec 2.0 □□□□□□□□□
2. K-means □ 128 □□□□□□□
3. □□□□□□ ID □□□□□□
4. □□□□□□□□□□□□

**□ □□□□□**

**□□□□□□□□**

- □□□□□□□MiniBatchKMeans
- □□□□□□n_clusters=128, random_state=42, batch_size=2048, n_init=3
- □□□□□(□□□□□□, 768) -□□□□□□□□□□□□
- □□□128 □□□□□□□□□□□□+□□□□
- □□□□□joblib.dump □□□□ pickle □□

**□□□□□□**

1. □□□□□□□: all_frames.shape = (15,234, 768)
2. KMeans □□: 128 □□□□□□□
3. □□□□□□□: cluster_centers_.shape = (128, 768)
4. □□□□:□□□□□→ □□□□□□□ ID (0-127)

> **□ □□**
>
> □□□cat□□□□□□□[ID:52, ID:23, ID:78]□□□□□□□□□□□□□□□□□□□□□□□[k æ t]□□□□□□□□□

**5.** □□□□□□□□□□

**RQ1**

☒☒☒☒☒☒☒☒VQ ☒☒☒☒☒☒☒☒☒☒☒☒☒☒☒ ☒☒☒☒☒☒☒☒☒☒☒☒☒☒☒☒☒☒☒☒☒

☒ ☒☒

☒☒☒☒☒☒☒☒☒☒☒☒☒☒☒☒☒☒☒☒☒☒☒☒☒☒☒☒☒☒☒☒☒☒☒☒☒☒☒☒☒☒☒☒☒☒☒☒☒☒☒☒☒☒☒☒☒☒☒☒

# 5.1 RQ1􏰀􏰁􏰂􏰃􏰄􏰅

􏰀􏰁􏰂􏰃

- 􏰀􏰁􏰂􏰃wav2vec2-base-960h 􏰄􏰅􏰆􏰇
- 􏰀􏰁􏰂􏰃􏰄􏰅 vs VQ 􏰆􏰇􏰈(vs 􏰉􏰊􏰋)
- 􏰀􏰁􏰂􏰃􏰄􏰅(􏰆􏰇􏰈􏰉􏰊􏰋􏰌)

􏰀􏰁􏰂􏰃

- **F1** 􏰀􏰁􏰂􏰃􏰄􏰅􏰆􏰇􏰈􏰉􏰊􏰋
- 􏰀􏰁􏰂􏰃􏰄􏰅
- 􏰀􏰁􏰂􏰃􏰄􏰅􏰆􏰇􏰈􏰉􏰊􏰋

**RQ2**

☒☒☒☒☒☒☒☒☒☒☒☒☒☒☒☒☒☒☒☒☒☒☒☒☒☒☒☒☒ ☒☒☒☒☒☒☒☒☒☒☒☒☒☒☒☒☒☒☒☒☒☒☒☒☒ (Panchendrarajan & Zubiaga, 2024)

# 6. ☒☒☒☒☒☒☒☒

> □ □□
>
> □□□□□□□□□□□□□Proof of Concept□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□
> □□□□□□□□□□□□□

## □□□□

- **Docker + Poetry□**□□□□□□□□□□□□□
- □□□□□□CPU □□□□MacBook Pro□- GPU □□□□□□□□□
- □□□□□□□□□□□□□□□□ 100 □□□□□□□□□□□□□□□□□□
- □□□□□□□□□□□3 □□□□□□□□□□□□□□□□□□

## □□□□□□

| □□□□□□ | □□□□□□ | □□ | □□ |
|---|---|---|---|
| **LibriSpeech** | 100 | □□□□□□□□□□ | □□□□□□□□ |

| □□□□□□ | □□□□□ | □□ | □□ |
|---|---|---|---|
| **Common Voice** | 100 | □□□□□□□□□□ | □□□□□□□□□□□ |

> **□ □□□□□**
>
> - LibriSpeech□□□□□□□□□□□□□□□□□□□□□□□
> - Common Voice□Mozilla □□□□□□□□□□□□□□□□□
>
> ⟹ □□□□□□□□□ 16kHz □□□□□□□□□□□□□

## RQ1□□□□□□□□□□□□

□ □□□□□

□□□□□□□□□

- G2P-EN □□□□□□□□text → phoneme list□
- □□□□□□□□□□□□□□np.linspace □□
- □□□□-□□□□□□□□□□□□□□

□□□□□□

- □□□□wav2vec2 □□□□□□768 □□□
- □□□□VQ □□□□ ID□0-127 □□□□
- □□□□□□□□□□□□□□□□□□□□□□□□

□□□□□□

- □□□□□□□□□□□□□□□□□□□□□□
- □□□train-test split (70%-30%)
- □□□63 □□□□□□□□□□□□

## RQ2□□□□□□□□□□□□□

□ □□□□□

□□□□□□□□□□□□□□□□□□

- □□□wav2vec2 □□□□□□□□□□□□768 □□□
- □□□□□StandardScaler □□
- □□□□□□□□□□□□8 □□□□□□□

□□□□□□□□□□□□+□□□□□□

- □□□□□□□□□□□□□□□768 □□□
- □□□□□□F0 □□□□□□□□□□□□□□
- □□□□□librosa.pyin □□□□□□□□□
- □□□□□□□□ 772 □□□□□□□□□□□
- □□□□□□□□□ StandardScaler □□

## □□□□ 1□□□□□□□□□□

> **□ □□□□□**
>
> **LibriSpeech test.clean □RQ1 □□□**
> - Hugging Face □□□□□□□□□□□□□□□□
> - □□□ 100 □□□□□□□□
> - □□□□□□□□□□□□□□□□□□□□□□□□□
> - □□□□□□□□□□+□□□□□□
>
> **Common Voice 13.0 □RQ2 □□□**
> - Mozilla □□□□□□□□□□□□□
> - □□□□□□□□□□□□□□□□□□□□
> - □□□□□□teens, twenties, thirties, forties, fifties, sixties, seventies, eighties
> - □□□□□□□□□□+□□□+□□□□□

⬚ ⬚⬚⬚⬚⬚

⬚⬚⬚⬚⬚⬚⬚

LibriSpeech⬚

```
{
  "file": "6930-75918-0000.flac",
  "audio": {"array": [-6.10e-05, 9.15e-05, ...], "sampling_rate": 16000},
  "text": "CONCORD RETURNED TO ITS PLACE AMIDST THE TENTS",
  "speaker_id": 6930
}
```

Common Voice⬚

```
{
  "audio": {"array": [0.001, -0.002, ...], "sampling_rate": 48000},
  "sentence": "The quick brown fox jumps over the lazy dog",
  "age": "twenties"
}
```

□□□□ 2□□□□□□

> □ □□□□□
>
> □□□□□□□
>
> 1. wav2vec2-base-960h □□□□□□□□
> 2. □□□□□□□16kHz □□□□□□□□□□□□□
> 3. □□□□□□□□Shape (□□□□□, 768 □□)
> 4. □□□librispeech_micro_continuous.npy
>
> **VQ □□□□□□**
>
> 1. □□□□□□□□□□□□Shape (□□□□□□, 768)
> 2. MiniBatchKMeans □□□128 □□□□□□
> 3. □□□□□□□□□vq_kmeans_128_micro.pkl
> 4. □□□□□□□□□□□□□□

□ □□□□□

**□□□□□□**

- 1 □□: (149, 768) → 149 □□□□×768 □□□□□□□
- 100 □□□□: (15,234, 768) → □ 15,234 □□□□
- VQ □□□□□□: (128, 768) → 128 □□□□□□□□

**VQ □□□□**

- □□□□□□: [0.73, −0.45, 1.23, …] (768 □□)
- → □□□□□ ID: 25 (0–127 □□□)
- VQ □□□□□□□□□□: [25, 25, 25, 52, 52, 78, 78, …]

**□□□□ 3□□□□□□□ notebook □□□□**

| Notebook | □□ | □□□□ |
|---|---|---|
| rq1_probing_pipeline | □□□□□□ | □□□ vs VQ □□□□□□□□□□□□ |

| Notebook | ☒☒ | ☒☒☒☒ |
|---|---|---|
| rq2_hybrid_model_poc | ☒☒☒☒☒☒☒☒ | ☒☒☒☒☒☒☒+☒☒☒☒☒☒☒☒ |

□□□□□□□□□

# 6.3 □□□□□□□□□□

□ □□□□□

□□□□□□□□□□□□□□□□□□

□□□□□□□□□□□□□□□□□□: /workspace/data/.cache
[□□□□] LibriSpeech□□□□□□□□ 100 □□□□□□□□□□□□□□□□
[□□□□] □□□LibriSpeech□□□□□□□□□:
{'file': '6930-75918-0000.flac', 'text': 'CONCORD RETURNED...'}
LibriSpeech□□□□□□ 100 □□ /workspace/data/raw/librispeech_micro □□□□□□□□□□□□

□□□□□□□□□□□□□

[□□□□] K-Means□□all_frames□shape: (15234, 768), Dtype: float32
VQ□□□□ (KMeans) □ 128 □□□□□□□□□□...
[□□□□] □□□□□□□□shape: (128, 768)
VQ□□□□ outputs/models/vq_kmeans_128_micro.pkl □□□□□□□□□□

□□□□□□□

- Poetry □□□□ Python □□□□□
- Docker □□□□□□□□ OS □□□□□□□
- requirements □□□□□□ □□□□□□□□
- Hugging Face datasets/transformers □□□□□□□□

# 7. □□□□□□□

□ □□

□□□□□□□□□□□□□□□

□□□□□□

- □□□□□63 □□□□□□□□□□□□
- □□□□□□□□LibriSpeech micro (100 □□□□□□□ 33,464 □□□□)
- □□□□□□□□□□G2P-EN +□□□□□□□□□□□□□□
- □□□□□train-test split (70%-30%) □□□□□□□□□□

**□□□□□□wav2vec 2.0□□□□□**

- □□□□□768 □□□□□□□□□□□□□□□□□□□
- □□□□□□□□□□□□□□□□□□□□□□□□□□
- □□□□□□□□□□□□□□□□□□□□□□□□
- □□□□□□□□□□□□1/63≈1.6%□□□□□□□□□□

□ □□□□□

□□□□□□□□□

- G2P-EN□□□□□ "A MAN SAID..." → □□□['AH0', ' ', 'M', 'AE1', 'N', ...]
- □□□□□np.linspace □□□□□□□□□□□□□□□□□
- □□150 □□□□□10 □□→ □□□ 15 □□□□□□□
- □□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□

□□□□□□□□□□□□□□□□□□□

- □□□□□□ "CONCORD RETURNED TO ITS PLACE AMIDST THE TENTS"
- □□□□□42 □□□['K', 'AA1', 'N', 'K', 'AO2', 'R', 'D', ' ', 'R', 'IH0', 'T', 'ER1', 'N', 'D', ...]
- □□□□□□□175 □□□□
- □□□□□[0, 4, 8, 12, 16, 20, 25, 29, ...] → □□□□□ 4 □□□□□□□□
- □□□□□□□□□□33,464 □□□□□□ 100 □□□× 768 □□□□□

# 7.1 RQ1：□□□□□□□□□□□□□□□□□□

## VQ □□□□□□□□□□□

□ □□

□□□□□□□□□□

- VQ □□□□□□128 □□□□□□□□□□□ 2 □□□□
- □□□□□□□□□□→ □□□□□□□□ ID（0-127 □□□□

**VQ □□□□□□□□□□□**

- □□□□□□□□□□□□□□□□□□□□□□□□□□□
- □□□□□□□□768 □□□→1 □□□□□□□□□□□□
- □□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□

**□ □□**

**VQ □□□□□□□□□□**

**□□□□□□□□**

- □□□□□□□□□□□□□□□□□□□□□□
- 768 □□→1 □□□□□□□□□□□□
- □□□□□□□□□□□□□□□□ ID =□□□□□□□□□□
- □□□□□□□□□□□□□□

**□□□ VQ □□□□**

- □□□□□: [0.73, −0.45, 1.23, …] (768 □□)
- VQ □□:□□□□ ID=52 (□□□□□□□)
- □□□cat□: [ID:52, ID:23, ID:78] ⟶ □□□□□□
- □□□ID=52→/k/□ID=23→/æ/□ID=78→/t/ □□□□□□

□ □□

□□□□□□□□□□□□□□□□□□□

□□□□□□

- □□□□□□□□□Common Voice micro (100 □□□□□□□□□□□□□)
- □□□□□8 □□□□□□□□□□(teens, twenties, thirties, forties, fifties, sixties, seventies, eighties)
- □□□□□□□□□□□□ vs □□□□□□□
- □□□□train-test split +□□□□□□□□□□+ classification_report

□ □□

□□□□□□□□□□□□□□□□□□□□□□□□

- □□□□Wav2Vec2 □□□□□□□□□□□□□□□□□□768 □□□
- □□□□StandardScaler □□□
- □□□□□8 □□□□□□□□□□□□□□□□
- □□□□□□□□□□□□□□□□□□□ 1/8 = 12.5%

□ □□

□□□□□□□□□□□□□□+□□□□□□□□□□□

- □□□□□□□□□□□□(768 □□) +□□□□□(4 □□) = 772 □□
- □□□□□□□□□□wav2vec2 □□□□□□□□□□768 □□□□
- □□□□□□F0 □□□□□□□□□□□□□□□□□□□□□□□□□□□□
- F0 □□□librosa.pyin(fmin=C2, fmax=C7) □□□□□□□□□□□□
- □□□□□□np.hstack □□□□□□□→ 772 □□□□□□□□
- □□□□□□□□□□ StandardScaler □□□

▨ ▨▨

▨▨▨▨▨▨▨▨

- X_neural: (100, 768) -▨▨▨▨▨▨▨▨▨
- X_acoustic: (100, 4) - [mean_f0, std_f0, jitter, shimmer]
- X_hybrid: (100, 772) -▨▨▨▨▨▨▨▨▨▨▨
- ▨▨▨▨▨▨▨≈0.0, ▨▨▨▨≈1.0 ▨▨▨▨▨▨▨▨

□□□□□□□□□□□□□□

> ☐ □□
>
> □□□□□□□□□□□
>
> - □□□□□□□□□□□□1/8 = 12.5% (8 □□□□□)
> - □□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□
> - □□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□
> - □□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□
>
> □□□□□2□□□□□□□□□
>
> - □□□□□□□□'twenties'
> - □□□□□□□□□(768,) □□□□□□□□□□
> - □□□□□[mean_f0: 192.33, std_f0: 15.7, jitter: 0.02, shimmer: 0.1]
> - □□□□□□(772,) =□□□□□□(768) +□□(4)

□ □□

□□□□□□□□□□□□□

**RQ1 □□□□□**

- □□□□□□□□□□□□□□wav2vec 2.0 □□□□□□□□□□□□□
- VQ □□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□
- □□□□□□□□□□→□□□□□□□□□□□□□□□□□□□□□□

**RQ2 □□□□□**

- □□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□
- □□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□
- □□□□□□□□□□□F0 □□□□□□□□□□□□□□□□□□

☒ ☒☒

☒☒☒☒☒☒☒☒☒☒☒☒☒☒

- ☒☒☒☒☒☒☒☒☒☒☒☒☒☒☒☒☒☒☒☒☒
- ☒☒☒☒☒☒☒☒☒☒☒☒☒☒☒☒☒☒☒☒☒☒☒
- ☒☒☒☒☒☒☒☒☒☒☒☒☒☒☒☒☒☒☒☒☒☒☒☒☒

□ □□

□□□□□□□□□□□□

1. □□□□□numpy □□□□□□□□□□□ 1 □□□□□□
2. □□□□□□□□torchaudio.transforms.Resample □□
3. □□□□Wav2Vec2Processor □□□□□□□□□□□□□
4. □□□□□torch.no_grad() □□ GPU □□□□□□
5. □□□□CPU □□□numpy □□□□□□□□□□

# 7.3 □□□□□□□

□□□□□

| □□ | □□ | □□□ |
|---|---|---|
| □□□□□□□□ | G2P-EN □□□□□□□□□ | Montreal Forced Aligner □□ |
| □□□□□□□ | 100 □□□□×2 □□□□□□ | □□□□□□□□□□□ |
| □□□□ | CPU □□ | □□□□ GPU □□ |
| □□□□□ | wav2vec2-base | WavLM-Large □□ |

> **□ □□□□□**
>
> **Montreal Forced Aligner (MFA)** =□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□ G2P-EN □□□□□□□□□□□□□□□□□□□□□□□□□

**8.** □□□□□□□□

## □□□□

- MFA □□□□□□□□□□□□□□□
- □□□□□□□□□□□□□□□
- WavLM-Large □□□□□□□
- □□□□□

## □□□□

- □□□□□□□□□□□□□□□□□□□□□□
- □□□□□□□□□□□□
- □□□□□□□□□□□□

# 8.2 □□□□□□□

## □□□□□□□

- SSL □□□□□□□□□□□□□□□□□□□□□
- □□□□□□□□□□□□□□□□□□□□□□□□ (Panchendrarajan & Zubiaga, 2024; Tsvilodub et al., 2025)
- □□□□□□□□□□□□□□□□□□□

## □□□□□□□□

- □□□□□□□□□□□□□□□□
- □□□□□□□□□□□□□□□□□□□
- □□□□□□□□□□□□□□ (Jarosz, 2019)

# 9. ░░░

□ □□

□□□□□□□□□□□□□□□□□

1. □□□□□□□□□□□□□□
2. □□□□□□□□□□□□□
3. □□□□□□□□□□□□□□□□□□□
4. □□□□□□□□□□□□□□□□□
5. □□□□□□□□□□□□□

□□□□

- □□□□□□□□□□□
- □□□□□□□□□
- □□□□□□□□

□ □□

□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□

1. □□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□
2. □□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□
3. □□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□
4. □□□□□□□□□□□□□□□□□□□□□□□□□□

□□□□□□□□□□□□□□□□□ □□□□□□□□□□□□□□□□□□□□□□

# 9.2 □□□□□□□

□□□□□□□□□□□□□□□□

□ □□□□□□

□□□□□□□□□□□

- data/processed/librispeech_micro_continuous.npy□□□□□□□□□□
- outputs/models/vq_kmeans_128_micro.pkl□□□□□□ VQ □□□□
- outputs/figures/cm_Continuous.png□□□□□□□□□□□□
- outputs/figures/cm_Discrete (VQ).png□VQ □□□□□□□□□□□□
- outputs/figures/cm_□□□.png, cm_□□□ (VQ).png□□□□□□□□

□□□□□□□□□□

- □□□□□□63 □□□□□'  ', 'AA0', 'AA1', 'AE1', 'AH0', ...□
- □□□□□□□□□□33,464 □□□□□×768 □□□□□□□□□33,464×1 □□□□□□□□
- □□/□□□□□□□23,424/10,040 □□□□□
- VQ □□□□□128 □□□□□□(128, 768)□□□□□□□□□□□

□□□□□□□

- □□□□ notebooks/prepare.ipynb □□□□□□□
- Docker □□□□□□□□□□□□□□□
- random_state □□□□□□□□□□□□□□□□□

Note: Reference file path needs to be adjusted for compilation

## ▨▨▨▨

Astrach, G., & Pinter, Y. (2025, June). *Probing Subphonemes in Morphology Models* (Issue arXiv:2505.11297). arXiv. https://doi.org/10.48550/arXiv.2505.11297

Baevski, A., Hsu, W.-N., Conneau, A., & Auli, M. (2022, May). *Unsupervised Speech Recognition* (Issue arXiv:2105.11084). arXiv. https://doi.org/10.48550/arXiv.2105.11084

Begŭs, G. (2020). Generative Adversarial Phonology: Modeling Unsupervised Phonetic and Phonological Learning With Neural Networks. *Frontiers in Artificial Intelligence*, *3*. https://doi.org/10.3389/frai.2020.00044

Chen, J., & Elsner, M. (2023, May). *Exploring How Generative Adversarial Networks Learn Phonological Representations* (Issue arXiv:2305.12501). arXiv. https://doi.org/10.48550/arXiv.2305.12501

Cho, C. J., Lee, N., Gupta, A., Agarwal, D., Chen, E., Black, A. W., & Anumanchipalli, G. K. (2025, March). *Sylber: Syllabic Embedding Representation of Speech from Raw Audio* (Issue arXiv:2410.07168). arXiv. https://doi.org/10.48550/arXiv.2410.07168

Choi, K., Pasad, A., Nakamura, T., Fukayama, S., Livescu, K., & Watanabe, S. (2024, June). *Self-Supervised Speech Representations Are More Phonetic than Semantic* (Issue arXiv:2406.08619). arXiv. https://doi.org/10.48550/arXiv.2406.08619

Gosztolya, G., Kiss-Vetráb, M., Svindt, V., Bóna, J., & Hoffmann, I. (2024). *Wav2vec 2.0 Embeddings Are No Swiss Army Knife-A Case Study for Multiple Sclerosis.*

Higy, B., Gelderloos, L., Alishahi, A., & Chrupała, G. (2021). Discrete Representations in Neural Models of Spoken Language. In J. Bastings, Y. Belinkov, E. Dupoux, M. Giulianelli, D. Hupkes, Y. Pinter, & H. Sajjad (Eds.), *Proceedings of the Fourth BlackboxNLP Workshop on Analyzing and Interpreting Neural Networks for NLP: Proceedings of the Fourth BlackboxNLP Workshop on Analyzing and Interpreting Neural Networks for NLP*. https://doi.org/10.18653/v1/2021.blackboxnlp-1.11

Jarosz, G. (2019). Computational Modeling of Phonological Learning. *Annual Review of Linguistics*, *5*(1), 67–90. https://doi.org/10.1146/annurev-linguistics-011718-011832

Kolachina, S., & Magyar, L. (2019). What Do Phone Embeddings Learn about Phonology?. In G. Nicolai & R. Cotterell (Eds.), *Proceedings of the 16th Workshop on Computational Research in Phonetics, Phonology, and Morphology: Proceedings of the 16th Workshop on Computational Research in Phonetics, Phonology, and Morphology*. https://doi.org/10.18653/v1/W19-4219

Medin, L. B., Pellegrini, T., & Gelin, L. (2024). Self-Supervised Models for Phoneme Recognition: Applications in Children's Speech for Reading Learning. *Interspeech 2024*, 5168–5172. https://doi.org/10.21437/Interspeech.2024-1095

Mohamed, A., Lee, H.-y., Borgholt, L., Havtorn, J. D., Edin, J., Igel, C., Kirchhoff, K., Li, S.-W., Livescu, K., Maaløe, L., Sainath, T. N., & Watanabe, S. (2022). Self-Supervised Speech Representation Learning: A Review. *IEEE Journal of Selected Topics in Signal Processing*, *16*(6), 1179–1210. https://doi.org/10.1109/JSTSP.2022.3207050

Panchendrarajan, R., & Zubiaga, A. (2024, March). *Synergizing Machine Learning & Symbolic Methods: A Survey on Hybrid Approaches to Natural Language Processing* (Issue arXiv:2401.11972). arXiv. https://doi.org/10.48550/arXiv.2401.11972

Pasad, A., Chien, C.-M., Settle, S., & Livescu, K. (2024). What Do Self-Supervised Speech Models Know About Words?. *Transactions of the Association for Computational Linguistics*, *12*, 372–391. https://doi.org/10.1162/tacl_a_00656

Pouw, C., Kloots, M. d. H., Alishahi, A., & Zuidema, W. (2024). Perception of Phonological Assimilation by Neural Speech Recognition Models. *Computational Linguistics*, *50*(3), 1557–1585. https://doi.org/10.1162/coli_a_00526

Silfverberg, M. P., Mao, L., & Hulden, M. (2018). Sound Analogies with Phoneme Embeddings. *Society for Computation in Linguistics*, *1*(1). https://doi.org/10.7275/R5NZ85VD

Staples, R., & Graves, W. W. (2020). Neural Components of Reading Revealed by Distributed and Symbolic Computational Models. *Neurobiology of Language (Cambridge, Mass.)*, *1*(4), 381–401. https://doi.org/10.1162/nol_a_00018

Tsvilodub, P., Hawkins, R. D., & Franke, M. (2025, June). *Integrating Neural and Symbolic Components in a Model of Pragmatic Question-Answering* (Issue arXiv:2506.01474). arXiv. https://doi.org/10.48550/arXiv.2506.01474

Venkateswaran, N., Tang, K., & Wayland, R. (2025, June). *Probing for Phonology in Self-Supervised Speech Representations: A Case Study on Accent Perception* (Issue arXiv:2506.17542). arXiv. https://doi.org/10.48550/arXiv.2506.17542