# Unsupervised Active Learning: Optimizing Labeling Cost-Effectiveness for Automatic Speech Recognition

*Zhisheng Zheng[1], Ziyang Ma[1], Yu Wang[2,3], Xie Chen[1*]*

[1]MoE Key Lab of Artificial Intelligence, AI Institute, X-LANCE Lab
Department of Computer Science and Engineering, Shanghai Jiao Tong University
[2]Cooperative Medianet Innovation Center, Shanghai Jiao Tong University
[3]Shanghai AI Laboratory

{zzs666, zym.22, yuwangsjtu, chenxie95}@sjtu.edu.cn

## Abstract

In recent years, speech-based self-supervised learning (SSL) has made significant progress in various tasks, including automatic speech recognition (ASR). An ASR model with decent performance can be realized by fine-tuning an SSL model with a small fraction of labeled data. Reducing the demand for labeled data is always of great practical value. In this paper, we further extend the use of SSL to cut down labeling costs with active learning. Three types of units on different granularities are derived from speech signals in an unsupervised way, and their effects are compared by applying a contrastive data selection method. The experimental results show that our proposed data selection framework can effectively improve the word error rate (WER) by more than 11% with the same amount of labeled data, or halve the labeling cost while maintaining the same WER, compared to random selection.

**Index Terms**: speech recognition, self-supervised learning, fine-tuning, unsupervised data selection

## 1. Introduction

Self-supervised learning (SSL) has emerged as a promising machine learning paradigm that allows us to learn more robust and distinctive features from unlabeled data, by leveraging inherent structure inside data without explicit labels. Recent works [1, 2, 3] have demonstrated that SSL models can extract high-quality and generalizable speech representations for various downstream speech-related tasks, such as speech recognition, speaker verification, and emotion recognition [4].

The SSL model typically consists of two-stage training: pre-training and fine-tuning. In the pre-training phase, several studies [5, 6] have demonstrated that using higher quality unlabeled speech data can improve the model's generalization performance. On the other hand, to achieve high performance in downstream tasks like ASR, it is necessary to fine-tune the pre-trained model with task-specific labeled data. However, acquiring labeled data in the fine-tuning stage can be expensive and challenging. Therefore, a practical challenge is how to select domain-relevant or task-relevant speech data within a limited budget for annotation in order to maximize the cost-effectiveness of labeling, making SSL models more practical and accessible.

In this paper, we present a completely unsupervised framework for selecting domain-relevant speech data. As illustrated in Fig. 1, the process involves generating discrete to-
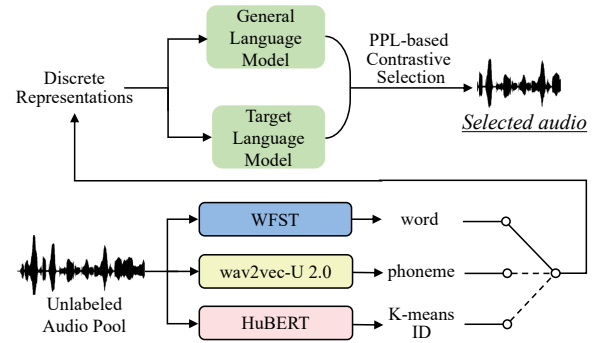
Figure 1: *The pipeline of unsupervised data selection on different granularities. The solid line represents we switch to word level. In this case the general language model is trained with word sequences from the WFST decoder, and the target language model is obtained by fine-tuning the general model with domain-specific text. Different granularities levels (K-means ID, phoneme, and word) are investigated.*

ken sequences at different levels of granularities (K-means ID, phoneme, word) from unlabeled speech data using intermediate models such as HuBERT, wav2vec-U 2.0, and WFST. We then calculate the perplexity (PPL) of these sequences using two pre-trained language models (a general LM and a target LM). Finally, we apply our PPL-based contrastive data selection approach to select the speech data that is most relevant to the target text.

Similar to previous works [5, 7], we use the discrete representations as the input units for our language model. Baevski et al. [8] pointed out that the discrete representations generated at different stages may capture different levels of acoustic information, which could potentially affect the accuracy of language modeling and the effectiveness of subsequent data selection. Therefore, in addition to the final recognition accuracy, other factors such as the amount of labeled data required for fine-tuning and the computational complexity of the process should also be considered when selecting the most appropriate granularities level for speech recognition.

We demonstrate the efficacy of our proposed unsupervised data selection method in the fine-tuning stage on a subset of GigaSpeech [9]. At the same level of granularity, by selecting only 100 hours of speech audio that closely match the given corpus, and fine-tuning the HuBERT base model on this labeled data without using any language models, we are able to reduce the Word Error Rate (WER) by more than 11% on all evaluated target domains. The main contributions can be summarized in three folds:

- We propose a novel, completely unsupervised active learning framework for speech data selection, which effectively reduces the cost of data labeling.
- We analyze the impact of different granularities levels on data selection and measure the trade-off between process complexity and recognition accuracy.
- Our proposed framework can either reduce the WER by over 11% with the same amount of labeled data, or cut the labeling cost to half while maintaining the same WER, compared to random selection.

# 2. Related work

## 2.1. Unsupervised data selection

Unsupervised data selection is a crucial technique aimed at achieving the goal of reducing the need for labeled data while maintaining high performance in downstream tasks for a specified target domain. In the field of natural language processing (NLP), various approaches [10, 11, 12] have been proposed for unsupervised data selection, including domain adaptation and topic models. In the field of ASR, however, the discrete representations of speech are not explicit, which reduces the possibility of borrowing methods from the NLP field. Therefore, a key challenge is how to obtain discrete token representations from continuous speech signals for speech data slection. Lu et al. [5] encoded speech signals into acoustically discrete tokens via self-supervised learning frameworks. Park et al. [6] calculated frame-level losses on a target data set and a training data set through SSL models, and then averaged these losses at the utterance level for subsequent selection. In addition to SSL-based methods, traditional unsupervised methods are still applicable. Drugman et al. [13] selectd data with low confidence scores from a speech recognition system. Malhotra et al. [14] proposed an entropy-based method for selecting the data that is most informative and uncertain for ASR.

## 2.2. SSL models

Self-supervised learning has attracted a lot of attention in recent years due to its potential to overcome the limitations of supervised learning that require large amounts of labeled data. It can be viewed as a two-stage process: pre-training and fine-tuning. In the pre-training stage, a model is trained on a large amount of unlabelled data using diverse self-supervised criteria, such as generative [15, 16, 17]; contrastive [2, 18, 19]; predictive [1, 3, 20, 21]. These tasks help the model learn general representations from the data without human labels. In the fine-tuning stage, the pre-trained model is continuous to be trained using a smaller amount of labeled data on the target domain and the representations are transferred to adapt a specific downstream task, ultimately leading to improved performance.

## 2.3. wav2vec-U 2.0

Wav2vec-U 2.0 [22] is an enhanced ASR system with a simplified architecture that achieves higher accuracy without requiring any pre-processing on the audio-side. Similar to its predecessor [8], wav2vec-U 2.0 learns the structure of speech from unlabeled audio data via self-supervised speech representations derived from either wav2vec 2.0 [2] or XLSR [23] models. These speech representations are then mapped to phonemes through a Generative Adversarial Network (GAN).
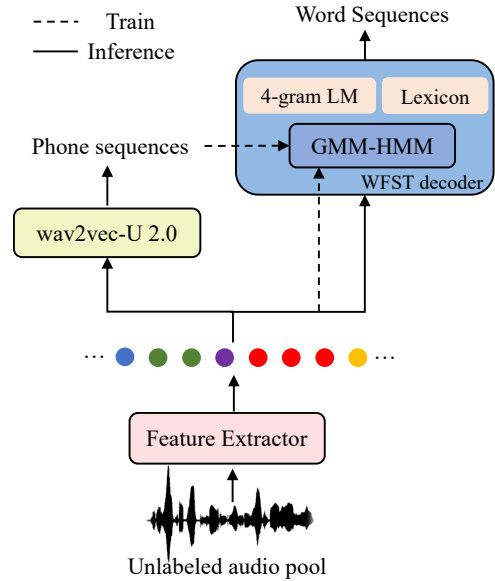
# 3. Methods



Figure 2: *The process of decoding audio into word sequences. Feature extractor is the pre-trained wav2Vec 2.0 Large (LV-60) model. 4-gram language model and lexicon are from the data preparation stage for training the wav2vec-U 2.0 model.*

In this section, we begin by addressing the process of obtaining phoneme labels from raw audio data in § 3.1, followed by a description of our transducer that decodes the unlabeled audio into words in § 3.2. Both of these sections are illustrated in Fig. 2. The unsupervised data selection strategy is elaborated upon in § 3.3.

## 3.1. Phoneme Recognizer

Wav2vec-U 2.0 is a highly effective unsupervised ASR system that stands out for its ability to take raw representations extracted from SSL models such as wav2vec 2.0 or HuBERT as input and output the corresponding sequence of phonemes for a given speech signal.

Prior to extracting features from speech using SSL models, it is necessary to perform VAD [24] on the speech signal to improve recognition accuracy. Subsequently, the SSL models like wav2vec 2.0 are used as feature extractors to obtain representations from the preprocessed speech. These extracted features are then fed into a pre-trained wav2vec-U 2.0 model, which outputs a sequence of phonemes as the final transcription.

## 3.2. HMM-based Transcription

Section § 3.1 has outlined the process of obtaining high-quality phoneme transcriptions, which can be used as targets for training a sufficiently robust GMM-HMM model. Instead of using Mel Frequency Cepstral Coefficient (MFCC) feature as the input to the GMM-HMM model, we use the same frame-level representations from the feature extractor as described in § 3.1.

In order to generate word sequences, we train a 4-gram language model and create a lexicon from a public text corpus, which is also used for building the text input for the wav2vec-U 2.0 model. These are then utilized to construct a Weighted Finite State Transducer (WFST) system in combination with the

previous trained HMM model. This system allows us to generate word pseudo-sequences from the unlabelled audio data.

### 3.3. Contrastive data selection

Contrastive data selection is a technique aimed at selecting samples from a larger dataset that are most relevant or similar to a target domain or task. Unlike [5], where acoustically discrete labels are used as LM input, we utilize sub-word units (BPE) to train two language models (LM). The first LM is trained on sub-word level pseudo labels since they are more relevant to our task. Although we attempt to use publicly available text datasets as the training corpus for the model, the results are mediocre. The second language model is obtained by simply fine-tuning the first LM using a limited sample of domain-specific data.

With these two trained LMs, we calculate the PPL of each sentence individually. However, our contrastive data selection algorithm does not directly use these sentence-level PPL, but instead uses the PPL of audio-level. Data selection at the utterance level may be more complex and inaccurate, and may require more prior knowledge because a single piece of text cannot fully represent the topic, sentiment, semantics, and other aspects of a speech, especially if the text is of low quality with high WER. Compared to that, audio-level selection method can mitigate the impact of irrelevant information, such as short sentences composed of common words and multi-topic words, thus allowing our algorithm to focus more on selecting topic-related audio.

The equation for the perplexity-based contrastive selection of each audio is defined as follows:

$$\eta = \frac{\overline{PPL_{LM2}} - \overline{PPL_{LM1}}}{\overline{PPL_{LM1}}}$$

where $\overline{PPL}$ denotes the average perplexity calculated for all utterances in every audio.

Subsequently, we select the audio within a budget based on the ascending order of $\eta$. See Fig. 1 for an illustration.

## 4. Experiments

### 4.1. Datasets

#### 4.1.1. LibriSpeech and GigaSpeech

The LibriSpeech corpus [25] is a widely-used speech dataset that contains approximately 1,000 hours of transcribed audio data from read English audio books. The GigaSpeech corpus [9] is a large-scale multi-domain English dataset that consists of over 10,000 hours of high quality labeled audios, covering a diverse topics, such as *Crime, Science, News*, etc.

#### 4.1.2. Cross-Domain Dataset

This paper presents a medium-sized dataset, consisting of a 1,000-hour cross-domain subset of GigaSpeech. The dataset is unique in its multi-source, multi-style composition, with each theme comprising an equal amount of data.

The dataset we have compiled includes 4 topics, namely *Crime, Health and Fitness, Howto and Style, and Science and Technology*. To ensure that the dataset is topic-balanced, we have included 100 hours of audio data for each of the four topics in the training set. To further augment the dataset and bring the total amount of data to 1000 hours, we have added an additional 600 hours of audio data from audiobooks, podcasts and youtube that are not specific to any of the four topics to the

training set. In addition to the training set, we have constructed dedicated validation and test sets for each of the 4 topics. Each of these sets contains 5 hours of audio data, providing ample material for model development, evaluation, and comparison.

For the integrity and representativeness of the training, validation, and test sets, we have taken care to avoid any overlap in audio segments between these sets during the sampling process. Specifically, for the validation and test sets, we have sampled from the *M*-size GigaSpeech training subset, as the word error rate of this subset is 0%. For the training set, however, we require 100 hours of topic-specific data for each of the four topics, which cannot be fully provided by the *M*-size and *L*-size subsets. Therefore, we have sampled the remaining audio data from the *XL*-size subset, resulting in approximately 270 hours of audio data with a word error rate of 4%.

### 4.2. Setup

For our experiments, we use the pre-trained wav2Vec 2.0 Large (LV-60) model as the feature extractor to obtain speech representations. We use these representations extracted from 960 hours of LibriSpeech and 30,000 randomly selected text samples from a publicly available training corpus to train a sufficiently good GAN model in the wav2vec-U 2.0 system [22]. The final model achieves a PER of 10.9% on the LibriSpeech dev-other. We use this model to decode the training set of our Cross-Domain dataset and obtain phoneme pseudo-sequences, which have a PER of approximately 18.8% (as shown in Fig. 3).

To improve recognition accuracy, we train a GMM-HMM model with speech representations extracted from the SSL model as input and utilize the pseudo-labels as targets, resulting in a final PER of 15.4%. We construct a WFST decoder by combining the GMM-HMM model, 4-gram language model, and lexicon. This allows us to decode the raw audio data to obtain word-level transcriptions, which results in a WER of around 32.4% (as shown in Fig. 3).
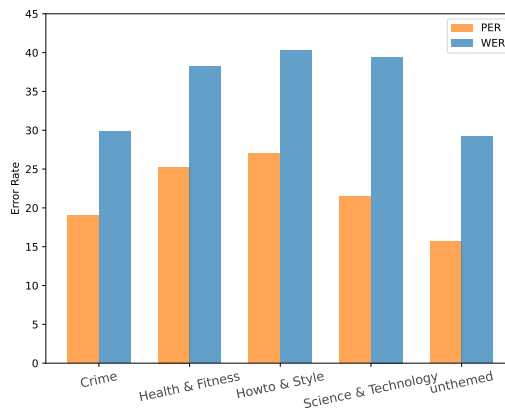


Figure 3: *Phone and word error rates on different categories. Phoneme recognition results are obtained from wav2vec-U 2.0 decoding, while word recognition results are from the decoding of the WFST decoder.*

As shown in Fig. 1, we extract the representations from the sixth layer [1] of the HuBERT base model to obtain K-Means IDs. We then apply the K-means algorithm to cluster the representations into 500 classes, and use the resulting K-means clustering IDs as direct input for the language models.

For phoneme-based discrete representations, we also directly use them as inputs for the language model without any further processing. However, considering the scalability of the word vocabulary and the performance of our small model, we constructed a sub-word corpus with a vocabulary size of 5000 using the BPE algorithm.

We then train our first language model using discrete token corpus, employing a Long Short-Term Memory (LSTM) [26] with 2 layers of hidden units and a vector dimension of 768. Several key hyperparameters of this model are set as follows: the learning rate is 1.0e-4, the number of epochs is set to 10 and the dropout rate is 0.2. As for the second language model, we simply fine-tune the first one using the given domain text with the same lexicon.

To evaluate the quality of our selected 100 hours of data, we employ the off-the-shelf pre-trained HuBERT base model as the quality assessor. We fine-tune each model for 80,000 steps and use the Viterbi algorithm as the decoding method.

### 4.3. Results

Table 1: *WERs of contrastive data selection on different granularities. All reported results are obtained by fine-tuning the HuBERT base model on 100 hours of labeled data and utilizing the Viterbi algorithm without language models.*

| Data Selection Algorithm | Crime | Health and Fitness | Howto and Style | Science and Technology |
|---|---|---|---|---|
| Random | 7.11 | 8.87 | 8.96 | 9.12 |
| Categorized[1] | 6.04 | 7.58 | 8.16 | 7.91 |
| **PPL-based Contrastive Selection (Granularities)** | | | | |
| K-means id | 6.40 | 7.73 | 8.13 | 8.22 |
| Phoneme | 7.39 | 8.85 | 8.3 | 8.63 |
| Words | 6.26 | 7.57 | 7.96 | 8.12 |
| Words*[2] | 6.05 | 7.65 | 7.97 | 8.01 |

[1] Data labeled with domain-specific classification tags.
[2] Ground truth word sequences.

Table 1 shows the WERs on the test set of our *Cross-Domain* dataset, with different data selection strategies and discrete token granularities levels. All results are obtained from decoding the test set without any language models, using the same HuBERT base model fine-tuned on 100-hour labeled data. The first two rows show the results of a random selection of labeled data and data labeled with domain-specific classification tags, respectively. The following rows show the results of the PPL-based contrastive selection method with different discrete token granularities levels, including K-means id, phoneme, and word-level tokens. The last row shows the results of using ground truth word sequences as LM input. In general, the PPL-based contrastive selection method outperforms the random sampling in almost all cases, regardless of granularities levels. The use of word-level tokens yields the best results (relatively more than 11%) across all domains, with ground truth word sequences performing even better, while phoneme-level tokens result in the highest WERs.

Table 2 compares the impact of labeled data of varying durations on the fine-tuning performance of SSL models, using the same word-level contrastive data selection method. Compared to a random sample of 100 hours of labeled data using our framework, we can achieve similar performance with only 50 hours of labeled data, which means the cost of labeling has

Table 2: *WERs with different amounts of labeled data. Results are evaluated for labeled data of different durations, using the same data selection algorithm.*

| Data Selection Algorithm | Labeled data | Crime | Health and Fitness | Howto and Style | Science and Technology |
|---|---|---|---|---|---|
| Random | 100h | 7.11 | 8.87 | 8.96 | 9.12 |
| Words | 100h | 6.26 | 7.57 | 7.96 | 8.17 |
| | 80h | 6.47 | 7.83 | 8.36 | 8.44 |
| | 60h | 6.83 | 8.24 | 8.85 | 9.02 |
| | 50h | 7.02 | 8.59 | 8.98 | 9.13 |

been cut in half.

### 4.4. Granularity Analysis

Based on the results of Table 1 and the difficulty in obtaining discrete representations, we can conduct a granularities analysis to determine which level of granularities is the most effective for speech recognition. Although using pseudo words as the granularities achieves the best performance and this performance can be further improved as the WER decreases, the process requires a multi-step inference procedure, which may be challenging to implement. In such situations, K-means ID-level representations may be a more practical alternative as they are the easiest to obtain. Compared with word sequences, HuBERT K-means IDs are derived from the acoustic features and may capture more detailed information about the acoustic characteristics of the speech signal. This may be the reason why using K-means IDs can achieve good performance.

## 5. Discussion

In this work, we made initial attempts to explore how to perform data selection based on completely unsupervised methods. Although the proposed algorithm has achieved good performance, there are several intriguing aspects that are worth investigating:

- Is there an optimal level of granularities for discrete tokens in speech data selection?
- Can this data selection algorithm be applied to other speech-related downstream tasks beyond ASR?
- Is it possible to develop a more optimal and simpler unsupervised data selection strategy for active learning?

We will research these issues in the future.

## 6. Conclusion

In recent years, self-supervised learning (SSL) for speech has demonstrated promising results in enhancing different downstream tasks such as speech recognition. In this paper, we investigate the problem of reducing the labeling cost while maintaining high performance in ASR through efficient data selection for SSL fine-tuning within a limited budget. We present a fully unsupervised and flexible active learning framework that selects relevant data based on the perplexity-based contrastive selection method. We analyze and compare the effectiveness of our framework using three different levels of granularities for discrete tokens: K-means ID, phoneme, and word. The optimal level is determined based on the selection performance and the complexity of the process. Our experimental results confirm the effectiveness of our framework for SSL fine-tuning data selection, which achieves significant improvements in WER while being more cost-effective in terms of annotation.

# 7. References

[1] W.-N. Hsu, B. Bolte, Y.-H. H. Tsai, K. Lakhotia, R. Salakhutdinov, and A. Mohamed, "HuBERT: Self-supervised speech representation learning by masked prediction of hidden units," in *Proc. ICASSP 2021*, 2021, pp. 3451–3460.

[2] A. Baevski, Y. Zhou, A. Mohamed, and M. Auli, "wav2vec 2.0: A framework for self-supervised learning of speech representations," *Proc. NIPS 2020*, pp. 12 449–12 460, 2020.

[3] A. Baevski, W.-N. Hsu, Q. Xu, A. Babu, J. Gu, and M. Auli, "Data2vec: A general framework for self-supervised learning in speech, vision and language," in *Proc. ICML 2022*, 2022, pp. 1298–1312.

[4] S. wen Yang, P.-H. Chi, Y.-S. Chuang, C.-I. J. Lai, K. Lakhotia, Y. Y. Lin, A. T. Liu, J. Shi, X. Chang, G.-T. Lin, T.-H. Huang, W.-C. Tseng, K. tik Lee, D.-R. Liu, Z. Huang, S. Dong, S.-W. Li, S. Watanabe, A. Mohamed, and H. yi Lee, "SUPERB: Speech Processing Universal PERformance Benchmark," in *Proc. Interspeech 2021*, 2021, pp. 1194–1198.

[5] Z. Lu, Y. Wang, Y. Zhang, W. Han, Z. Chen, and P. Haghani, "Unsupervised data selection via discrete speech representation for ASR," in *Proc. Interspeech 2022*, 2022, pp. 3393–3397.

[6] C. Park, R. Ahmad, and T. Hain, "Unsupervised data selection for speech recognition with contrastive loss ratios," in *Proc. ICASSP 2022*, 2022, pp. 8587–8591.

[7] C.-C. Chiu, J. Qin, Y. Zhang, J. Yu, and Y. Wu, "Self-supervised learning with random-projection quantizer for speech recognition," in *Proc. ICML 2022*, 2022, pp. 3915–3924.

[8] A. Baevski, W.-N. Hsu, A. Conneau, and M. Auli, "Unsupervised speech recognition," *Proc. NIPS 2021*, pp. 27 826–27 839, 2021.

[9] G. Chen, S. Chai, G.-B. Wang, J. Du, W.-Q. Zhang, C. Weng, D. Su, D. Povey, J. Trmal, J. Zhang, M. Jin, S. Khudanpur, S. Watanabe, S. Zhao, W. Zou, X. Li, X. Yao, Y. Wang, Z. You, and Z. Yan, "GigaSpeech: An evolving, multi-Domain ASR corpus with 10,000 hours of transcribed audio," in *Proc. Interspeech 2021*, 2021, pp. 3670–3674.

[10] A. Ramponi and B. Plank, "Neural unsupervised domain adaptation in NLP—A survey," in *Proc. ICCL 2020*, 2020, pp. 6838–6855.

[11] D. M. Blei, A. Y. Ng, and M. I. Jordan, "Latent Dirichlet allocation," *Proc. JMLR 2003*, pp. 993–1022, 2003.

[12] H. Jelodar, Y. Wang, C. Yuan, X. Feng, X. Jiang, Y. Li, and L. Zhao, "Latent Dirichlet allocation (LDA) and topic modeling: models, applications, a survey," *Proc. MTAs 2019*, pp. 15 169–15 211, 2019.

[13] T. Drugman, J. Pylkkonen, and R. Kneser, "Active and semi-supervised learning in ASR: Benefits on the acoustic and language models," *arXiv preprint arXiv:1903.02852*, 2019.

[14] K. Malhotra, S. Bansal, and S. Ganapathy, "Active learning methods for low resource end-to-end speech recognition." in *Proc. Interspeech 2019*, 2019, pp. 2215–2219.

[15] X. Yue and H. Li, "Phonetically motivated self-supervised speech representation learning." in *Proc. Interspeech 2021*, 2021, pp. 746–750.

[16] L. Liu and Y. Huang, "Masked pre-trained encoder base on joint CTC-Transformer," *arXiv preprint arXiv:2005.11978*, 2020.

[17] A. H. Liu, Y.-A. Chung, and J. Glass, "Non-autoregressive predictive coding for learning speech representations from local dependencies," in *Proc. Interspeech 2021*, 2021, pp. 3730–3734.

[18] A. v. d. Oord, Y. Li, and O. Vinyals, "Representation learning with contrastive predictive coding," *arXiv preprint arXiv:1807.03748*, 2018.

[19] Y.-A. Chung, Y. Zhang, W. Han, C.-C. Chiu, J. Qin, R. Pang, and Y. Wu, "W2v-bert: Combining contrastive learning and masked language modeling for self-supervised speech pre-training," in *Proc. ASRU 2021*, 2021, pp. 244–250.

[20] S. Chen, C. Wang, Z. Chen, Y. Wu, S. Liu, Z. Chen, J. Li, N. Kanda, T. Yoshioka, X. Xiao *et al.*, "Wavlm: Large-scale self-supervised pre-training for full stack speech processing," *Proc. IEEE J Sel Top Signal Process 2022*, pp. 1505–1518, 2022.

[21] Z. Ma, Z. Zheng, C. Tang, Y. Wang, and X. Chen, "MT4SSL: Boosting self-supervised speech representation learning by integrating multiple targets," *arXiv preprint arXiv:2211.07321*, 2022.

[22] A. H. Liu, W.-N. Hsu, M. Auli, and A. Baevski, "Towards end-to-end unsupervised speech recognition," in *Proc. SLT 2022*, 2022, pp. 221–228.

[23] Q. Xu, A. Baevski, and M. Auli, "Simple and effective zero-shot cross-lingual phoneme recognition," *arXiv preprint arXiv:2109.11680*, 2021.

[24] Z.-H. Tan, A. Kr. Sarkar, and N. Dehak, "rVAD: An unsupervised segment-based robust voice activity detection method," *Proc. Computer speech & language 2020*, pp. 1–21, 2020.

[25] V. Panayotov, G. Chen, D. Povey, and S. Khudanpur, "Librispeech: An ASR corpus based on public domain audio books," in *Proc. ICASSP 2015*, 2015, pp. 5206–5210.

[26] S. Hochreiter and J. Schmidhuber, "Long short-term memory," *Proc. Neural computation 1997*, pp. 1735–1780, 1997.