



# On Combining Global and Localized Self-Supervised Models of Speech

Sri Harsha Dumpala<sup>1,2</sup>, Chandramouli S Sastry<sup>1,2</sup>, Rudolf Uher<sup>2,3</sup> and Sageev Oore<sup>1,2</sup>

<sup>1</sup>Vector Institute, <sup>2</sup>Dalhousie University and <sup>3</sup>Nova Scotia Health, Canada

{sriharsha.d, cssastry, uher, sageev}@dal.ca

## Abstract

Self supervised learning involves learning general-purpose representations that can be useful in a variety of downstream tasks. In this work, we study the application of speech-embeddings derived from popular self-supervised learning frameworks such as wav2vec-2.0 and HuBERT over four different speech-classification tasks such as sentiment classification, command detection, emotion classification and depression detection. We distinguish between and discuss self-supervised training tasks that induce *localized* and *global* features of speech based on their temporal granularity: noting that self-supervised representation learning frameworks based on the masked language-modeling objective – such as wav2vec-2.0 and HuBERT – induce *localized* embeddings, we define a self-supervised learning framework based on SimSiam for learning global features of speech. Through our evaluations, we find that these global representations are better suited for tasks such as depression detection and emotion classification while the *localized* embeddings of speech can be very useful in tasks such as speech-command detection; we also find that our proposed model outperforms TRILL – a popular model for learning global representations. Finally, we also propose and confirm empirically that combining the global and localized representations of speech helps obtain better performance across a range of downstream tasks than each of the individual embedding methods.

**Index Terms:** Simple Siamese (SimSiam), self-supervised learning, emotion classification, sentiment classification, depression detection

## 1. Introduction

Self-supervised learning frameworks aim to learn general purpose representations from large amounts of data through strategically defined training tasks that do not require per-instance manual annotations. Popular self-supervised learning frameworks include: contrastive learning frameworks [1, 2, 3, 4], clustering frameworks [5, 6] and non-contrastive learning frameworks [7, 8]. In this work we will explore the use of simple siamese representation learning (SimSiam) [8], a non-contrastive representation learning framework, for speech.

As part of the self-supervised learning setup, models are trained on cleverly defined pretext tasks to obtain general purpose representations which can then be used for the downstream tasks. In the context of learning self-supervised representations of speech, we can categorize the pretext tasks as:

- **Localized** pretext tasks focus on learning segment-level representations of speech conditional on remaining contextual segments of the speech. As an example, self-supervised learning frameworks over raw speech waves that use masked prediction as the pretext task [9, 10, 11] are said to induce localized representations of speech as the task involves learning to represent and identify segments of speech conditional on the remaining segments;

these localized representations of speech can be used to obtain state-of-the-art performance in automatic speech recognition (ASR) [12, 9, 10]. In summary, examples of localized pretext tasks for speech include: predicting the content of the unseen regions [13, 14], language model-style pre-training [15], predicting discrete targets of masked regions [16, 9, 10], combining self-supervised learning with adversarial training [17].

- **Global** pretext tasks focus on learning global representations of entire speech wave. For example, [18] proposes to contrastively learn representations of mel-spectrograms derived from speech segments – on the other hand, wav2vec-2.0 contrastively learns representations of speech segments conditional on remaining segments of the raw speech. In a similar vein, [19] proposes a self-supervised learning strategy based on BYOL [7] for learning representations of spectrograms derived from speech segments without requiring the use of negative samples (i.e. without contrastive learning). Some other examples of self-supervised learning frameworks based on global pretext tasks [20, 11, 21, 22] are usually 1) applied to mel-spectrograms; and 2) use contrastive learning frameworks which rely on negative samples, thus requiring large batch sizes to train the models.

Intuitively, localized features of speech effectively represent fine details of speech segments conditional on remaining speech segments – the state-of-the-art performance in ASR obtained by the use of localized features support this intuition. On the other hand, global features of speech represent more temporally stable elements of speech and as we show in our results, these global representations of speech are better suited than localized representations to downstream tasks such as emotion classification and depression detection. In fact, TRILL[11] was proposed to learn such non-semantic representations of speech. Accordingly, the localized representations of speech depend on the speech duration and change quickly throughout a speech sample, whereas global representations do not. In general, a global representation of speech can be obtained from the localized features of speech – such as those obtained from BERT-style training – by averaging across the constituent speech segments. We note that we can draw a parallel to text-representation learning: for example, BERT training over text corpora learns context-sensitive embeddings of words whereas some downstream tasks require sentence embeddings; see [23] for a novel fine-tuning method over BERT for deriving sentence embeddings from pre-trained BERT.

In this paper, we will propose SimSiam-speech, a simple siamese self-supervised learning framework for speech to learn global features of speech. We show that these global features achieve better performance than localized features on some downstream tasks. We will also show that combining these global features with localized features will achieve improved performance on various downstream tasks.

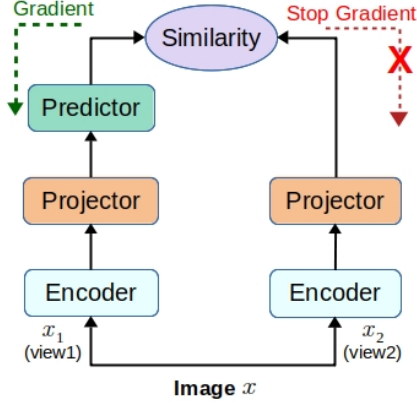


Figure 1: Outline of the SimSiam framework initially proposed for images [8]. Blocks with the same color have weight sharing.

The main contributions of this work are as follows:

- We introduce a SimSiam-based framework for speech (SimSiam-Speech): Most previous self-supervised learning work on speech has been based on frameworks such as masked prediction similar to BERT [24], simple contrastive learning [1] and bootstrap your own latent [7]. In this work, we have extended the SimSiam framework for speech representation learning.
- Building on the intuition that global and localized features of speech have different temporal granularity, we propose to combine these two representations to obtain better performance on downstream tasks. In our evaluations, we report on and analyse the application of localized and global – both, independently and together – self supervised representations of speech on a diverse set of speech classification tasks.

## 2. Background

The SimSiam framework [8] is an extension of BYOL [7] for learning self-supervised representations of images. Like other recently proposed self-supervised learning frameworks, SimSiam induces representations invariant over the set of predefined data augmentations. The core components of the SimSiam framework are as shown in Figure 1: given an input image  $x$ , the two random augmented views –  $x_1$  and  $x_2$  – of the input  $x$  are first passed through a series of Encoder-Projector-Predictor, and Encoder-Projector networks respectively; the training objective involves maximizing the cosine similarity between these final representations. The predictor and the stop-gradient are essential to avoid model collapse, which we will also show in the ablation studies. Finally, the encoder representations are used as the general-purpose representations of images in downstream tasks.

In summary, SimSiam is a simplified self-supervised learning framework which has been shown to achieve SOTA performance on image downstream tasks without relying on negative sample pairs, large batch sizes and momentum encoder. In this paper, we adapt the SimSiam framework to learn self-supervised global representations of speech.

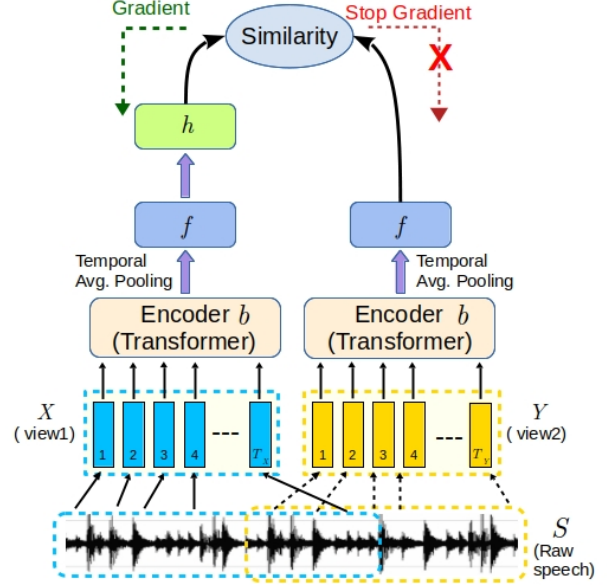


Figure 2: Outline of the proposed SimSiam framework for speech. Blocks with the same color have weight sharing.

## 3. Proposed Approach

**SimSiam Framework for Speech** builds on top of SimSiam (Figure 1) and is shown in Figure 2: given a raw speech waveform  $S$  as input, two random augmented sequences of speech segments –  $X$  and  $Y$  – are used to enforce invariance across a sequence of speech augmentations described below. We obtain the representations of speech by average pooling the encoder outputs across the time steps; accordingly, these representations are used as input to the Projector  $f$ . We refer to these representations as global representations of speech and apply them in downstream speech classification tasks.

Given a raw speech waveform  $S$  with sampling frequency  $f_s$ , we obtain a sequence of segments  $s_i \in \mathbb{R}^d$  such that  $s_i$  spans  $d/f_s$  seconds – in other words, each segment  $s_i$  contains  $d$  waveform amplitudes and the waveform  $S$  containing  $D$  amplitude samples gets partitioned into  $T = D/d$  segments. We generate the two augmented views by first partitioning the sequence of segments  $\{s_i\}_{i \in [1, T]}$  into two contiguous sequences  $X = \{x_i\}_{i \in [1, T_X]}$  and  $Y = \{y_i\}_{i \in [1, T_Y]}$  with an overlap of 50-80% segments as shown in Figure 2. The duration of the two augmented views  $X$  and  $Y$  is in the range of 2.0-2.5 seconds, and  $T_X$  and  $T_Y$  can be different. We can now define the outputs obtained after applying the encoder  $b$  (i.e. the transformer backbone), projector  $f$  and predictor  $h$  as  $z_X = h(f(\hat{b}(X)))$  and  $z_Y = \text{stopgrad}\{f(\hat{b}(Y))\}$  where  $\hat{b}$  refers to the average-pooled outputs obtained after applying the encoder  $b$ . The training objective involves maximizing the cosine similarity between  $z_X$  and  $z_Y$  as shown:

$$s(X, Y) = \frac{z_X}{\|z_X\|_2} \cdot \frac{z_Y}{\|z_Y\|_2}, \quad (1)$$

where  $\|\cdot\|_2$  refers to  $l_2$ -norm. Finally, to have a symmetric loss as defined in [8], the loss function is defined as:

$$L = s(X, Y) + s(Y, X). \quad (2)$$

**Augmentations:** For each of the sub-sequences  $X$  and  $Y$ , we apply the following sequence of augmentations:

1. **Gaussian Noise:** To each segment, apply additive white Gaussian noise (AWGN) with SNR between 0 – 10 db and then normalize the amplitude of the resulting segment using min-max normalization.
2. **Shuffling:** Randomly select and shuffle 20 – 40% of the segments.
3. **Masking:** Randomly select 20 – 40% of the segments and replace the segments with an  $\mathbb{R}^d$  vector sampled uniformly from the hypercube  $[0.9, 1.1]^d$ .
4. **Add silence segments:** Select  $N$  random locations and insert silence segments whose  $\mathbb{R}^d$  vector has all components set equal to the lowest amplitude of the original speech waveform  $S$ . We choose  $N$  to be equal to 10% of the sequence length (i.e.  $T_X$  or  $T_Y$ ).

**Encoder network  $b$ :** We use a multi-layer Transformer architecture [24] as the encoder network. The transformer model consists of 12 transformer-encoder layers, each containing 12 attention heads. Each segment (i.e.  $x_i$  and  $y_i$ ) is projected to a 768-dimensional vector before passing as input to the transformer; for example, the input and output of the transformer for sequence  $X$  is of shape  $T_X \times 768$ . The non-linear projection networks used within the transformer consist of one hidden layer having 2048 units with ReLU as the non-linearity. Finally, we use the standard self-attention [25] as the Multi-Head-Attention module with layer-norm normalization.

**Projector ( $f$ ) and Predictor ( $h$ ) networks:** We follow the architecture of the projector and predictor as defined in [8]. The projector network takes a 768-dimensional global representation as input and has 2048-dimensional output vector while the prediction network has both input and output dimensions equal to 2048: the projector has two hidden layers, each having 2048 ReLU units and batch normalization (BN) [26] in between; the prediction network has one hidden layer with 512 ReLU units and BN.

**Pre-training dataset:** For a fair comparison with the previous self-supervised learning models, our SimSiam-speech model is pre-trained using the LibriSpeech dataset [27]. LibriSpeech dataset is derived from public domain audio books which contains 960 hours of speech data collected from 2338 different speakers (1128 female and 1210 male speakers). In this work, we discard utterances with duration less than 3 seconds as we select two speech segments each of duration ranging from 2.0-2.5 seconds, with an overlap of 50-80% in duration, to train our models.

**Pre-training Details:** Input speech is sampled at  $f_s = 16\text{kHz}$ . We take raw speech waveform as input and use  $d = 1000$  for creating segments. Our models are pre-trained with a batch size of 480 using Adam optimizer with a learning rate of  $3e^{-4}$  and a weight decay of  $5e^{-5}$ . Furthermore, the learning rate has a cosine decay schedule [28].

## 4. Experiments

### 4.1. Downstream Tasks

We compare and evaluate the representations extracted from the pre-trained SimSiam-speech model on a diverse set of downstream tasks to test different aspects of the representations. The downstream tasks are summarized in Table 1 and are as follows:

Table 1: *Details of the datasets used for downstream tasks. #Class, #Sample and Avg. Len(s) refer to the number of classes, number of samples and average length in seconds, respectively*

Dataset	Task	#Class	#Sample	Avg. Len(s)
MOSI [29]	Sentiment	2	2199	4.2
MOSEI [30]	Sentiment	2	23453	7.28
Speech-Commands [31]	Command	12	100,503	1.0
IEMOCAP [32]	Emotion	4	3846	4.5
DAIC [33]	Depression	2	219	960

1. **Sentiment classification:** We use CMU Multimodal Opinion-Level Sentiment Intensity (CMU-MOSI) [29] and (2) CMU Multimodal Opinion Sentiment and Emotion Intensity (CMU-MOSEI) [30] datasets to evaluate the pre-trained models for the task of sentiment classification.
2. **Low-vocabulary speech recognition:** We use speech commands [31] dataset which contains short spoken commands for the task of low-vocabulary speech recognition.
3. **Speech emotion recognition:** We use IEMOCAP [32] dataset for the downstream task of speech emotion recognition.
4. **Depression detection:** DAIC-WOZ [33] dataset is used for the downstream task of depression detection.

### 4.2. Models for Comparison

We compare the performance of the pre-trained SimSiam-speech model with various publicly available SOTA speech representation models:

- **Mockingjay [14]:** Mockingjay, a transformer-based model, was trained using BERT-like loss to predict masked frames.
- **VQ-wav2vec [16]:** VQ-wav2vec uses online k-means clustering to quantize the dense representations for training a deep convolutional network with a BERT-like loss.
- **Wav2vec-2.0 [9]:** Wav2vec-2.0, a transformer-based model, was trained to predict the discrete representation of the masked frames using a BERT-like loss.
- **HuBERT [10]:** HuBERT, an extension of wav2vec-2.0, utilizes an offline clustering step to provide aligned target labels to train transformer model using BERT-like loss.
- **TRILL [11]:** TRILL, a transformer-based model, was trained on mel spectrograms using a triplet loss.

Note that VQ-wav2vec, wav2vec-2.0 and HuBERT takes raw speech waveform as input whereas Mockingjay and TRILL takes mel spectrogram as input. While Mockingjay, VQ-wav2vec, wav2vec-2.0 and HuBERT learn localized features, TRILL learns global features. Further, we also use the openS-MILE [34] and COVAREP [35] features as a baseline.

### 4.3. Evaluation Details

For each downstream dataset, we freeze the weights of the pre-trained models, and extract representations from the encoder for

Table 2: Performance (in terms of accuracy) on different downstream tasks by using the representations obtained from different pre-trained models. Details of tasks: MOSI and MOSEI for sentiment classification, Speech commands for command detection, IEMOCAP for emotion classification, and DIAC-WoZ for depression detection

Model	MOSI	MOSEI	Speech		
			Commands	IEMOCAP	DAIC-WOZ
COVAREP	46.6	52.9	48.2	48.2	53.3
OpenSMILE	51.3	58.1	39.8	51.6	56.8
Mockingjay	56.3	69.1	87.4	52.4	58.7
VQ-wav2vec	54.1	68.2	86.7	51.6	60.4
Wav2vec-2.0	60.8	70.4	89.6	55.6	61.9
HuBERT	61.9	71.1	90.2	54.9	61.3
TRILL	56.8	63.5	78.5	61.7	64.1
SimSiam-S (Ours)	59.0	65.8	85.1	62.2	66.4
Wav2vec-2.0 + HuBERT	62.3	71.3	90.5	56.5	62.2
Wav2vec-2.0 + TRILL	62.2	70.7	90.3	62.8	68.2
HuBERT + TRILL	63.7	71.4	90.9	62.6	67.8
TRILL + SimSiam-S	59.3	65.5	85.3	62.4	66.9
Wav2vec-2.0 + SimSiam-S	64.8	71.2	91.1	<b>63.7</b>	<b>71.2</b>
HuBERT + SimSiam-S	<b>65.9</b>	<b>71.8</b>	<b>91.4</b>	63.5	70.3

the downstream task. We apply a non-linear projection (i.e. linear projection followed by non-linearity) to these representations and then pass it as input to the softmax classifier. The number of hidden units are selected based on the performance on the validation set. For the case of combining the representations extracted from different pre-trained models, we simply concatenate the non-linear projections applied to the representations before passing it as input to the softmax layer.

#### 4.4. Results

We provide results obtained by performing 5-fold cross validation on each dataset, unless otherwise mentioned.

Table 2 shows the performance of the proposed SimSiam-Speech (SimSiam-S) compared to previous self-supervised models across four diverse set of downstream tasks. For the downstream tasks of emotion recognition and depression detection, proposed SimSiam-S and TRILL performs better than the other pre-trained models, with SimSiam-S outperforming TRILL. For the downstream tasks of sentiment classification and low-vocabulary speech recognition, HuBERT and wav2vec-2.0 performs better. This shows that the models learning global representations (SimSiam-S and TRILL) perform better on some tasks while models learning localized representations (HuBERT and wav2vec-2.0) perform better on other tasks. We can also observe in Table 2 that combining global and localized representations (HuBERT + SimSiam-S and wav2vec-2.0 + SimSiam-S) achieve significant improvements in performance compared to the case of combining two global (wav2vec-2.0 + HuBERT) or two localized representations (TRILL + SimSiam-S). For all the downstream tasks, combining SimSiam-S with local representations (HuBERT or wav2vec-2.0) achieves the best performance.

#### 4.5. Ablation Studies

We conducted ablation studies to analyze the importance of stop-gradient, projector and predictor networks when SimSiam framework is extended to speech—which is a sequential data. We trained different SimSiam-Speech models by discarding one component at a time. All the models were pre-trained using the 960 hours of LibriSpeech dataset. It can be observed from Ta-

ble 3 that all the three components i.e., projector, predictor and stop-gradient are required to avoid the model collapse. Even discarding a single component will lead to the model collapse. These observations are in agreement with those reported in [8].

Table 3: Ablation study results. Proj., Pred., SG refer to projector network, predictor network and stop-gradient, respectively. Results obtained for different SimSiam-Speech models trained by discarding one of the component

Proj.	Pred.	SG	Speech	
			Commands	IEMOCAP
✗	✓	✓	6.8	19.1
✓	✗	✓	7.3	18.5
✓	✓	✗	6.5	20.8
✓	✓	✓	85.1	62.2

## 5. Conclusion

In this work, we discuss the application of self-supervised embeddings to speech classification tasks as the downstream tasks. In doing so, we highlight the distinction between self-supervised learning frameworks that favor *global* versus *localized* embeddings of speech: in this study, pretext tasks which operate on smaller speech units are said to yield localized embeddings, while those operating on entire speech sequences are said to yield global embeddings. We define a self-supervised learning framework for learning global features of speech and propose combining both global and localized features of speech for obtaining more robust downstream performance. Our empirical experiments confirm that global embeddings of speech capture certain paralinguistic elements and can be particularly useful in tasks such as emotion classification and depression detection, while localized embeddings of speech can be effective in tasks such as sentiment classification and speech-command detection. Importantly, we find that *combining* the global and localized representation of speech leads to better performance across a range of tasks than each of the individual embedding methods alone.

## 6. References

- [1] T. Chen, S. Kornblith, M. Norouzi, and G. Hinton, "A simple framework for contrastive learning of visual representations," in *International conference on machine learning*. PMLR, 2020, pp. 1597–1607.
- [2] X. Chen, H. Fan, R. Girshick, and K. He, "Improved baselines with momentum contrastive learning," *arXiv preprint arXiv:2003.04297*, 2020.
- [3] K. He, H. Fan, Y. Wu, S. Xie, and R. Girshick, "Momentum contrast for unsupervised visual representation learning," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2020, pp. 9729–9738.
- [4] O. Henaff, "Data-efficient image recognition with contrastive predictive coding," in *International Conference on Machine Learning*. PMLR, 2020, pp. 4182–4192.
- [5] Y. Asano, C. Rupprecht, and A. Vedaldi, "Self-labelling via simultaneous clustering and representation learning," in *International Conference on Learning Representations*, 2019.
- [6] M. Caron, I. Misra, J. Mairal, P. Goyal, P. Bojanowski, and A. Joulin, "Unsupervised learning of visual features by contrasting cluster assignments," in *Thirty-fourth Conference on Neural Information Processing Systems (NeurIPS)*, 2020.
- [7] J.-B. Grill, F. Strub, F. Althé, C. Tallec, P. H. Richemond, E. Buchatskaya, C. Doersch, B. A. Pires, Z. D. Guo, M. G. Azar *et al.*, "Bootstrap your own latent: A new approach to self-supervised learning," *arXiv preprint arXiv:2006.07733*, 2020.
- [8] X. Chen and K. He, "Exploring simple siamese representation learning," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2021, pp. 15 750–15 758.
- [9] A. Baevski, Y. Zhou, A. Mohamed, and M. Auli, "wav2vec 2.0: A framework for self-supervised learning of speech representations," *Advances in Neural Information Processing Systems*, vol. 34, 2020.
- [10] 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," *IEEE/ACM Transactions on Audio, Speech, and Language Processing*, vol. 29, pp. 3451–3460, 2021.
- [11] J. Shor, A. Jansen, R. Maor, O. Lang, O. Tuval, F. de Chaumont Quiry, M. Tagliasacchi, I. Shavitt, D. Emanuel, and Y. Haviv, "Towards learning a universal non-semantic representation of speech," *Proc. Interspeech 2020*, pp. 140–144, 2020.
- [12] A. Van den Oord, Y. Li, and O. Vinyals, "Representation learning with contrastive predictive coding," in *32nd Conference on Neural Information Processing Systems (NeurIPS 2018)*. NeurIPS, 2018.
- [13] Y.-A. Chung and J. Glass, "Improved speech representations with multi-target autoregressive predictive coding," in *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, 2020, pp. 2353–2358.
- [14] A. T. Liu, S.-w. Yang, P.-H. Chi, P.-c. Hsu, and H.-y. Lee, "Mockingjay: Unsupervised speech representation learning with deep bidirectional transformer encoders," in *ICASSP 2020-2020 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 2020, pp. 6419–6423.
- [15] Y.-A. Chung, W.-N. Hsu, H. Tang, and J. R. Glass, "An unsupervised autoregressive model for speech representation learning," in *INTERSPEECH*, 2019.
- [16] A. Baevski, S. Schneider, and M. Auli, "vq-wav2vec: Self-supervised learning of discrete speech representations," in *International Conference on Learning Representations*, 2019.
- [17] A. Baevski, W.-N. Hsu, A. Conneau, and M. Auli, "Unsupervised speech recognition," *Advances in Neural Information Processing Systems*, vol. 34, 2021.
- [18] A. Saeed, D. Grangier, and N. Zeghidour, "Contrastive learning of general-purpose audio representations," in *ICASSP 2021-2021 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 2021, pp. 3875–3879.
- [19] D. Niizumi, D. Takeuchi, Y. Ohishi, N. Harada, and K. Kashino, "Byol for audio: Self-supervised learning for general-purpose audio representation," in *2021 International Joint Conference on Neural Networks (IJCNN)*. IEEE, 2021, pp. 1–8.
- [20] S. Chen, Y. Wu, C. Wang, Z. Chen, Z. Chen, S. Liu, J. Wu, Y. Qian, F. Wei, J. Li *et al.*, "Unispeech-sat: Universal speech representation learning with speaker aware pre-training," *arXiv preprint arXiv:2110.05752*, 2021.
- [21] D. Jiang, W. Li, M. Cao, W. Zou, and X. Li, "Speech simclr: Combining contrastive and reconstruction objective for self-supervised speech representation learning," *arXiv preprint arXiv:2010.13991*, 2020.
- [22] J. Shor, A. Jansen, W. Han, D. Park, and Y. Zhang, "Universal paralinguistic speech representations using self-supervised conformers," *arXiv preprint arXiv:2110.04621*, 2021.
- [23] N. Reimers and I. Gurevych, "Sentence-bert: Sentence embeddings using siamese bert-networks," in *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing*. Association for Computational Linguistics, 11 2019. [Online]. Available: <https://arxiv.org/abs/1908.10084>
- [24] J. D. M.-W. C. Kenton and L. K. Toutanova, "Bert: Pre-training of deep bidirectional transformers for language understanding," in *Proceedings of NAACL-HLT*, 2019, pp. 4171–4186.
- [25] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, Ł. Kaiser, and I. Polosukhin, "Attention is all you need," *Advances in neural information processing systems*, vol. 30, 2017.
- [26] S. Ioffe and C. Szegedy, "Batch normalization: Accelerating deep network training by reducing internal covariate shift," in *International conference on machine learning*. PMLR, 2015, pp. 448–456.
- [27] V. Panayotov, G. Chen, D. Povey, and S. Khudanpur, "Librispeech: an asr corpus based on public domain audio books," in *2015 IEEE international conference on acoustics, speech and signal processing (ICASSP)*. IEEE, 2015, pp. 5206–5210.
- [28] T. Chen, S. Kornblith, M. Norouzi, and G. Hinton, "A simple framework for contrastive learning of visual representations," in *International conference on machine learning*. PMLR, 2020, pp. 1597–1607.
- [29] A. Zadeh, R. Zellers, E. Pincus, and L.-P. Morency, "Multimodal sentiment intensity analysis in videos: Facial gestures and verbal messages," *IEEE Intelligent Systems*, vol. 31, pp. 82–88, 2016.
- [30] A. B. Zadeh, P. P. Liang, S. Poria, E. Cambria, and L.-P. Morency, "Multimodal language analysis in the wild: Cmu-mosei dataset and interpretable dynamic fusion graph," in *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, 2018, pp. 2236–2246.
- [31] P. Warden, "Speech commands: A dataset for limited-vocabulary speech recognition," *arXiv preprint arXiv:1804.03209*, 2018.
- [32] C. Busso, M. Bulut, C.-C. Lee, A. Kazemzadeh, E. Mower, S. Kim, J. N. Chang, S. Lee, and S. S. Narayanan, "Iemocap: Interactive emotional dyadic motion capture database," *Language resources and evaluation*, vol. 42, no. 4, pp. 335–359, 2008.
- [33] J. Gratch, R. Artstein, G. Lucas, G. Stratou, S. Scherer, A. Nazarian, R. Wood, J. Boberg, D. DeVault, S. Marsella *et al.*, "The distress analysis interview corpus of human and computer interviews," in *Proceedings of the Ninth International Conference on Language Resources and Evaluation (LREC'14)*, 2014, pp. 3123–3128.
- [34] F. Eyben, M. Wöllmer, and B. Schuller, "Opensmile: the munich versatile and fast open-source audio feature extractor," in *Proc. ACM conference on Multimedia*, 2010, pp. 1459–1462.
- [35] G. Degottex, J. Kane, T. Drugman, T. Raitio, and S. Scherer, "Covarep—a collaborative voice analysis repository for speech technologies," in *2014 IEEE international conference on acoustics, speech and signal processing (icassp)*. IEEE, 2014, pp. 960–964.