# Non-Contrastive Self-Supervised Learning of Utterance-Level Speech Representations

*Jaejin Cho[1], Raghavendra Pappagari[1], Piotr Żelasko[1,2], Laureano Moro-Velazquez[1],*
*Jesús Villalba[1,2], Najim Dehak[1,2]*

[1]Center for Language and Speech Processing, Johns Hopkins University, Baltimore, MD, USA
[2]Human Language Technology Center of Excellence, Johns Hopkins University, Baltimore, MD, USA

{jcho52,rpappag1,pzelasko,laureano,jvilla17,ndehak3}@jhu.edu

## Abstract

Considering the abundance of unlabeled speech data and the high labeling costs, unsupervised learning methods can be essential for better system development. One of the most successful methods is contrastive self-supervised methods, which require negative sampling: sampling alternative samples to contrast with the current sample (anchor). However, it is hard to ensure if all the negative samples belong to classes different from the anchor class without labels. This paper applies a non-contrastive self-supervised learning method on an unlabeled speech corpus to learn utterance-level embeddings. We used DIstillation with NO labels (DINO), proposed in computer vision, and adapted it to the speech domain. Unlike the contrastive methods, DINO does not require negative sampling. These embeddings were evaluated on speaker verification and emotion recognition. In speaker verification, the unsupervised DINO embedding with cosine scoring provided 4.38% EER on the VoxCeleb1 test trial. This outperforms the best contrastive self-supervised method by 40% relative in EER. An iterative pseudo-labeling training pipeline, not requiring speaker labels, further improved the EER to 1.89%. In emotion recognition, the DINO embedding performed 60.87, 79.21, and 56.98% in micro-f1 score on IEMOCAP, Crema-D, and MSP-Podcast, respectively. The results imply the generality of the DINO embedding to different speech applications.

**Index Terms**: self-supervised learning, speaker verification, emotion recognition, distillation, non-contrastive

## 1. Introduction

Self-supervised learning (SSL) is gaining more attention in many machine learning areas such as computer vision, natural language processing, and speech processing. SSL does not require labeled data for model training. In many works, fine-tuned/post-processed SSL models have shown promising results outperforming supervised methods when the same amount of labeled data is used [1, 2, 3, 4, 5, 6].

In speaker verification, different self-supervised methods have been proposed as in [7, 8, 9, 10, 11, 12, 13]. Some of these methods use a generative approach [7, 10, 8, 9], i.e., they learn to reconstruct the signal acoustic features from some latent representations. Usually, the goal is to factorize the information into frame-/segment-level and utterance-level latent factors, expecting that the former will encode phonetic information while the latter will encode the speaker information. For example, the work in [9] used an architecture based on Tacotron 2 multi-speaker text-to-speech to learn speaker embeddings.

SSL methods based on contrastive loss are also popu-
lar [11, 12, 13] in speaker verification. Contrastive losses intend to make the current sample (anchor) close to the augmented version of the anchor (positive sample) while making the positive sample farther from the negative samples in their embedding space. In this context, negative samples are desired to be different semantically from the positive sample. Since the samples are unlabeled, most contrastive SSL works in speaker verification compose negative samples just by randomly picking different samples to the anchor. However, in this random sampling, we are not sure if all the negative samples are from different classes w.r.t the positive sample. For example, when the anchor, thus also the positive sample, is an utterance from speaker A, there is a chance that some of the negative utterances come from speaker A as well. This could adversely affect the model training since the contrastive loss pushes the positive sample and negative sample farther to each other in the embedding space.

Non-contrastive methods, however, do not require negative samples, so they are free from the issue above. Moreover, non-contrastive methods have shown comparable or better performance compared to contrastive methods [5, 6]. Considering these, we propose to apply a non-contrastive SSL method originally proposed for computer vision (CV), DIstillation with NO labels (DINO) [6], that outperformed the previous SSL methods in many CV tasks, including linear and k-NN classification.

Since DINO training does not use explicit labels such as speaker or language IDs, we hypothesized the learned embedding to be general utterance-level embedding. In other words, the embedding may include attributes that are consistent within the utterance, such as speaker information, accent/language, emotion, and age. Thus, we evaluated the embedding not only for speaker verification but also for emotion recognition. The results confirmed that the DINO embedding includes both speaker and emotion information.

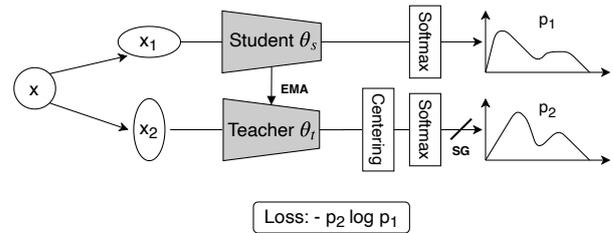## 2. Distillation with NO labels in speech



Figure 1: *DINO diagram. EMA stands for exponential moving average. SG stands for stop gradient. The figure shows a single augmented pair ($x_1$,$x_2$) for simplicity.*

In [6], the authors proposed a design to maximize the similarity between feature distributions of differently augmented images from an original image. This assumes that augmented images from one image keep the same semantic information. For example, although you crop two different portions from a dog image and make one a black and white image while making the other jittered, they are still dog images. This principle can be similarly applied to speech data. Assuming each utterance consists of a speech from one speaker, different segments extracted from the same utterance followed by noise augmentation share the same attributes that are consistent within the utterance, such as speaker information, accent/language, emotion, and age.

DINO follows a knowledge distillation scheme, where the outputs of a teacher network are used as ground truth to optimize a student network. However, typical knowledge distillation uses a pre-trained teacher network, while DINO trains both networks in parallel. Fig. 1 depicts DINO's training scheme. First, we augment a given utterance into a set of differently augmented segments, $S$. With speech data, for example, each segment is extracted randomly from the given utterance as either a short or long segment. Then, a sound such as babbling, music, noise in the background, or room impulse response effect is applied to each segment differently. The set $S$ includes two long segments, $s_1^l$ and $s_2^l$. All the augmented segments go through the student branch, while only the long segments go through the teacher branch. Each branch embeds information in the corresponding segments into the embedding vectors. Additionally, each network has a head with a softmax layer that classify each segment into a set of hypothetical classes. Thus, the student network $\theta_s$ is trained by minimizing,

$$\theta_s = \min_{\theta_s} \sum_{s \in \{s_1^l, s_2^l\}} \sum_{\substack{s' \in S \\ s' \neq s}} H(p_2(s|\theta_t), p_1(s'|\theta_s)). \quad (1)$$

where $H(a, b) = -a \cdot \log b$ is cross-entropy, and $p_1(.)$ and $p_2(.)$ are the softmax outputs of the student and teacher networks respectively. This loss intends to make the embeddings for all augmented versions of the utterance close between them, which relies on two assumptions. First, the long segments, used as input to the teacher, produce better representations than the short-segments. Second, the teacher network is always better than the student during the training, as explained below.

The neural network architecture for student and teacher models comprises a backbone followed by a projection head. The backbone can be any encoder that converts a sequence of vectors into a fixed-dimensional vector, e.g., using a pooling layer. The projection head consists of fully connected layers with non-linear activations. The student and teacher networks are initialized with the same architecture having the same parameters while they are updated in different ways during training. The student network is updated by gradient descent while the teacher network is updated by an exponential moving average on the student parameters, i.e., $\theta_t \leftarrow \lambda\theta_t + (1 - \lambda)\theta_s$. Parameter averaging is known to produce a better model [14, 15], and this is also the case with the teacher network to be better than the student network during the training. The student model aims at the distribution from the teacher network to improve. To avoid a model to find trivial solutions, i.e., having distributions where one dimension is dominant or having uniform distributions, *centering* and *sharpening* are applied. *centering* prevents one dimension from dominating by calculating a center and subtracting it from the logit before the softmax. However, *centering* encourages a uniform distribution, and thus

*sharpening* is also applied where it encourages peaky distributions. This is done by setting a low value for the temperature in the teacher softmax normalization.

# 3. Experimental setup[1]

## 3.1. Encoder pre-training using DINO

We used a light ResNet34 (LResNet34) encoder from [16] as the DINO backbone considering available resources, with minor modifications: a kernel size of 3 in the first convolution layer, a mean and standard deviation pooling, and a following affine layer that outputs a 256-dimensional vector used as a embedding vector. The classification head consists of 3 linear layers with their hidden dimension as 2048, followed by $l_2$ normalization and a weight normalized linear layer with 65536 softmax output dimension, which were the best setup in the original DINO paper [6] that we found work well with speech data, too.

For the augmentation, we first extracted 4 short and 2 long segments randomly from a given utterance where we set 2 and 4 s for short and long segment length, respectively. We chose the specific numbers for the extracted segments considering our available computational resource and training time, but extracting more short segments could improve performance, as in [6]. The extracted segments were augmented with babbling, music, noise in the background, and/or room impulse response effects.

VoxCeleb2 *dev* set of 5,994 speakers was used for training DINO without speaker labels. VoxCeleb2 corpus [17] consists of conversational speech utterances with moderate noise, which were processed from interview videos of 6,112 celebrities uploaded on Youtube, and it covers diverse ethnicity.

The final learned encoder was used for speaker verification and emotion recognition tasks employing other datasets.

## 3.2. Speaker verification

In the speaker verification experiments, we followed a pipeline to build a progressively improving speaker verification system, which does not require speaker labels [18]. Starting from the initial DINO model trained in section 3.1, the system development goes through the iterative clustering stage followed by the robust training stage. The training data was fixed as VoxCeleb2 *dev* without speaker labels for the whole pipeline.

In the iterative clustering stage, we trained a new larger model, ResNet34 [19] x-vector model, in a supervised way with the additive angular margin (AAM) [20] loss based on pseudo speaker labels generated using the initial DINO model. In detail, we extracted speaker embeddings from the initial model. Then, the embeddings were clustered using k-means clustering with 50k means, followed by agglomerative hierarchical clustering (AHC) with the number of clusters as 7500, which was heuristically determined for VoxCeleb2 *dev* [18]. The k-means clustering was used to make AHC computationally viable. Indices of the clusters are used as pseudo speaker labels for the supervised x-vector model training. The whole process was repeated 3 times until the speaker verification performance converged. During the process, the model parameters were continuously updated with the refined pseudo labels in each cycle. A 2-second segment was extracted per utterance to be used as a training sample.

In the robust training stage, we used a new larger model, Res2Net50 [21] with 26 for the width of filters (in the first residual bottleneck block) and 4 for the scale hyper-parameter. The model was trained in a supervised way with pseudo labels gen-

---

[1]The code for experiments will be uploaded to the public repository https://github.com/hyperion-ml/hyperion.

erated from the ResNet34 model after 3 cycles of the iterative clustering stage. After the first 30 epochs of training, the post-pooling layers of the model were fine-tuned with a larger margin, 0.5, in the AAM loss. A 2-second segment was extracted from each utterance to be used as a training sample, while 3-second segments were used in the large margin fine-tuning.

The learned embedding in each stage was evaluated for speaker verification on the original VoxCeleb1 *test* (voxceleb1_test_o), VoxSRC-21 *val*, or VoxSRC-21 *test* trials. The latter two trials are from VoxCeleb Speaker Recognition Challenge2021(VoxSRC-21), where the challenge has a special focus on multi-lingual verification. Our team's submission of the system having 6.88 in EER(%) in Table 2 ranked third in track 3: self-supervised speaker verification where the participants were allowed to develop systems only with VoxCeleb2 *dev* without using speaker labels.

### 3.3. Emotion recognition

The DINO model training uses segments (2s for $x_1$ and 4s for $x_2$ in Fig. 1) drawn from VoxCeleb2 *dev* dataset. For the model to preserve emotion information, $x_1$ and $x_2$ have to share the same emotion at least for the majority of training. As the median duration of utterances in VoxCeleb2 *dev* is only 6.08s, we expect consistent emotion between $x_1$ and $x_2$ in general. This assumption is also supported by the considered emotion datasets, which contain utterances of duration 2-11s with utterance-level emotion annotations. Hence, we hypothesize that the model preserves emotion-related information along with speaker identity.

To probe the DINO embeddings for emotion information, we evaluated them for the emotion recognition task. Specifically, we extracted DINO embeddings for three emotion recognition datasets, Crema-D [22], IEMOCAP [23], and MSP-Podcast [2] [24], and performed emotion classification using logistic regression on each corpus.

The Crema-D dataset consists of acted emotions with utterances ranging from 2-4s in duration; The IEMOCAP dataset is composed of utterances with induced emotions and mostly 2-7s in duration; The MSP-Podcast dataset contains spontaneous utterances of duration 3-11s. Regarding experimental setup, we performed leave-one-session-out cross-validation (5-fold CV) for IEMOCAP as in the previous works. For Crema-D, we used the same train/dev/test splits as in [25], and the standard splits for MSP-Podcast as in Release 1.4. We used *angry, happy, sad*, and *neutral* emotions in the IEMOCAP and Crema-D dataset, and an additional emotion class *disgust* in MSP-Podcast as in [25]. We report weighted average values of class-wise f1-scores – micro-f1 metric – for these experiments. For IEMOCAP, we report the average across five folds.

## 4. Results

### 4.1. Speaker verification

#### 4.1.1. DINO embedding in initial model training

In this experiment, we first compared two SSL embeddings: DINO and MoCo, in the initial model training stage in sec-

Table 1: *Comparison between DINO, momentum contrast (MoCo), and x-vector embeddings for speaker verification. The results are on the original VoxCeleb1 test with equal error rate (EER)(%) and MinDCF with $P_T$=0.01. The PLDA back-end was trained with VoxCeleb1 dev where its data size is 1/7 of VoxCeleb2 dev.*

|  | Cosine scoring | | PLDA | |
|---|---|---|---|---|
|  | EER(%) | MinDCF | EER(%) | MinDCF |
| DINO | 4.83 | 0.463 | 2.38 | 0.289 |
| MoCo [18] | 7.3 | - | - | - |
| x-vector | 1.94 | 0.207 | 1.88 | 0.189 |

tion 3.2. This stage does not include iterative clustering and pseudo labeling for supervised training. As shown in Table 1, DINO (non-contrastive method) outperforms MoCo (contrastive method) by 40% relatively in EER(%). This implies that an SSL model can learn embedding for speaker verification without negative sampling.

Next, we compared the performance of DINO to supervised x-vector in the DINO and x-vector rows of Table 1. Their encoder architectures were identical (LResNet34). Although x-vector performed better than DINO, it used VoxCeleb2-dev speaker labels while DINO with cosine scoring back-end did not use any labels at all. We also evaluated a PLDA back-end trained on VoxCeleb1-dev, which 7 times smaller than VoxCeleb. With PLDA where the gap between DINO and x-vector reduced further, while DINO only used 1/8 of the labels than x-vector. The data used for PLDA also can be employed to fine-tune the DINO encoder, expecting further improvement as it is a common observation in most of the SSL papers.

#### 4.1.2. Iterative clustering and robust training

The experimental results are shown in Table 2 along the process from the DINO initial model training to the robust training. In iterative clustering, the speaker verification performance is saturated around the 3rd iteration. This number of iterations until convergence is less than the one reported in [18][3], possibly due to starting from a better initial model. Thus, we generated the pseudo labels from the model after the 3rd iteration (ResNet34 (iter3)) to use them for the last model training in the robust training stage, which improved further to 1.89 in EER(%) on voxceleb1_test_o. This speaker verification system did not use any speaker labels in the development, and to compare, the supervised counterpart trained on about 2600 hours of speaker labeled data showed 0.93 in EER(%). Finally, the large-margin fine-tuning did not improve on the voxceleb1_test_o trial list, while it improved on VoxSRC-21 *val* and *test* trial lists.

### 4.2. Emotion recognition

Table 3 presents the results for emotion recognition using the DINO embedding. We compare our results with [25], where the authors trained a logistic regression classifier on top of pre-trained x-vectors for emotion recognition. Their x-vector model was pre-trained with 8kHz data to discriminate speakers. This includes telephone data other than VoxCeleb downsampled to 8kHz where the final number of speakers was around 13k [26]. We observed that DINO embeddings performed better than the x-vectors suggesting that DINO embeddings contain

---

[3]This is a rough comparison since detailed configurations are slightly different.

Table 2: *Speaker verification results over 3 different trial lists with progressing/different systems over the three stages. The numbers from [18] seems rounded to the nearest tenth. Pseudo labels for robust training were generated from ResNet34 (iter3).*

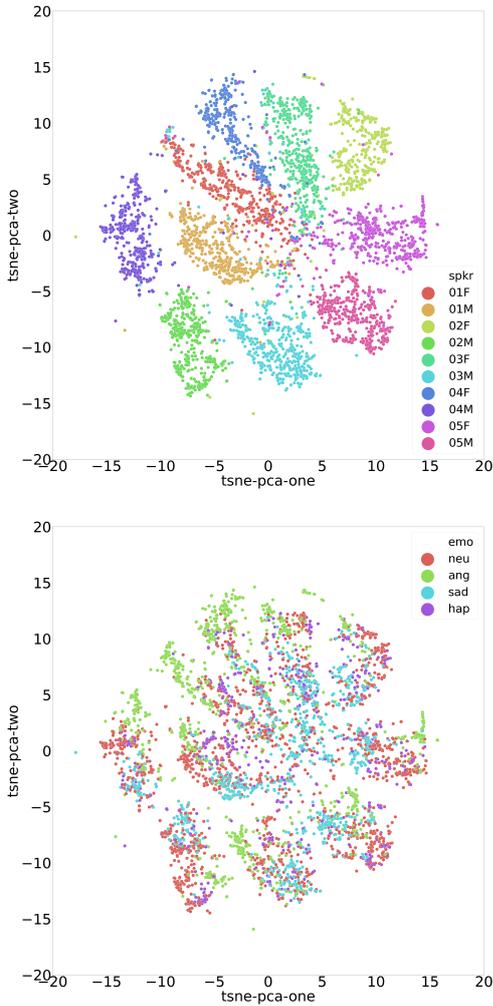| Stage | Algorithm/Loss | Model | EER (%) with cosine scoring | | |
|---|---|---|---|---|---|
| | | | voxceleb1_test_o | VoxSRC-21 *val* | VoxSRC-21 *test* |
| Initial model training (self-supervised learning) | DINO | LResNet34 | 4.83 | 13.96 | - |
| | MoCo | ECAPA [18] | 7.3 | - | - |
| Iterative clustering | AAM loss (margin=0.3) | ResNet34 (iter1) | 2.56 | 8.59 | - |
| | | ResNet34 (iter2) | 2.13 | 7.35 | - |
| | | ResNet34 (iter3) | 2.13 | 6.97 | - |
| | | ResNet34 (iter4) | 2.14 | 6.88 | - |
| | | ECAPA (iter7) [18] | 2.1 | - | - |
| Robust training + larg-margin fine-tuning | AAM loss (margin=0.5) | Res2Net50 | 1.89 | 6.50 | 6.88 |
| | | | 1.91 | 6.32 | 6.64 |



Figure 2: *Analysis of DINO embedding space for IEMOCAP using t-SNE plots. Each color represents one speaker in the top plot and one emotion in the bottom plot.*

Table 3: *Emotion classification results on three different dataset. All numbers in this table are micro-f1 (%) scores*

| | IEMOCAP | Crema-D | MSP-Podcast |
|---|---|---|---|
| x-vector [25] | 56.11 | 75.65 | 52.58 |
| DINO | **60.87** | **79.21** | **56.98** |

more emotion predictive information than the x-vectors. This result makes sense since supervised x-vectors are trained to retain only speaker information and also trained with more data. In contrast, the DINO model is trained to capture common information across extracted segments from the utterance.

Fig. 2 shows t-SNE plots of the DINO embedding space for the IEMOCAP dataset. From the top plot, we can observe clear clusters of speakers suggesting abundant speaker-relevant information in DINO embeddings. In the bottom plot, there are signs of emotion clusters for some speakers, especially for angry and sad, suggesting some emotional information is captured in the representations. Looking more closely within each speaker cluster, *angry* and *sad* are well separated which could be because they usually have distinct arousal levels (high for *angry* and low for *sad*).

## 5. Conclusion

In this paper, we learned utterance-level embeddings based on DINO, a non-contrastive self-supervised learning method originally proposed in CV. The embedding was evaluated on the speaker verification and emotion recognition tasks to check its generalizability. DINO embedding achieved the state-of-the-art result in speaker verification when no speaker labels were used, outperforming the previous best contrastive SSL embedding based on MoCo. Also, it reduced the number of iterations until convergence in the iterative clustering stage. When the DINO embedding was used for emotion recognition, it performed better than the x-vector embedding that was found to contain emotion-related information in a previous study. One thing to note is that the DINO embedding was learned without any labels while it still performed competitively with the embedding using speaker labels for both speaker verification and emotion recognition. The DINO method opens the way for leveraging unlabeled speech data, which is more easily available than labeled one.

Considering the DINO embedding may also embed other attributes consistent within a given utterance, we will test the embedding on other speech applications such as accent/language, speech pathology, and age recognition. Also, fine-tuning the DINO embedding to one of the applications with the labeled data is a natural expansion of this work.

## 6. Acknowledgements

# 7. References

[1] J. Devlin, M.-W. Chang, K. Lee, and K. Toutanova, "BERT: Pre-training of deep bidirectional transformers for language understanding," in *NAACL*, Jun. 2019, pp. 4171–4186. [Online]. Available: https://aclanthology.org/N19-1423

[2] J. Lee, W. Yoon, S. Kim, D. Kim, S. Kim, C. H. So, and J. Kang, "Biobert: a pre-trained biomedical language representation model for biomedical text mining," *Bioinformatics*, vol. 36, no. 4, pp. 1234–1240, 2020.

[3] A. Baevski, Y. Zhou, A. Mohamed, and M. Auli, "wav2vec 2.0: A framework for self-supervised learning of speech representations," in *NeurIPS*, vol. 33. Curran Associates, Inc., 2020, pp. 12 449–12 460. [Online]. Available: https://proceedings.neurips.cc/paper/2020/file/92d1e1eb1cd6f9fba3227870bb6d7f07-Paper.pdf

[4] T. Chen, S. Kornblith, K. Swersky, M. Norouzi, and G. E. Hinton, "Big self-supervised models are strong semi-supervised learners," *NeurIPS*, vol. 33, pp. 22 243–22 255, 2020.

[5] J.-B. Grill, F. Strub, F. Altché, C. Tallec, P. Richemond, E. Buchatskaya, C. Doersch, B. Avila Pires, Z. Guo, M. Gheshlaghi Azar, B. Piot, k. kavukcuoglu, R. Munos, and M. Valko, "Bootstrap your own latent - a new approach to self-supervised learning," in *NeurIPS*, vol. 33, 2020, pp. 21 271–21 284.

[6] M. Caron, H. Touvron, I. Misra, H. Jégou, J. Mairal, P. Bojanowski, and A. Joulin, "Emerging properties in self-supervised vision transformers," in *ICCV*, 2021.

[7] W.-N. Hsu, Y. Zhang, and J. Glass, "Unsupervised learning of disentangled and interpretable representations from sequential data," in *NeurIPS*, 2017, pp. 1876–1887.

[8] T. Stafylakis, A. J. Rohdin, O. Plchot, P. Mizera, and L. Burget, "Self-supervised speaker embeddings," in *Interspeech*, vol. 2019, no. 9, 2019, pp. 2863–2867. [Online]. Available: https://www.fit.vut.cz/research/publication/12092

[9] J. Cho, P. Żelasko, J. Villalba, S. Watanabe, and N. Dehak, "Learning Speaker Embedding from Text-to-Speech," in *Interspeech*, 2020, pp. 3256–3260.

[10] Z. Peng, S. Feng, and T. Lee, "Mixture factorized auto-encoder for unsupervised hierarchical deep factorization of speech signal," in *ICASSP*, 2020, pp. 6774–6778.

[11] J. Huh, H. S. Heo, J. Kang, S. Watanabe, and J. S. Chung, "Augmentation adversarial training for unsupervised speaker recognition," in *NeurIPS Workshop*, 2020.

[12] W. Xia, C. Zhang, C. Weng, M. Yu, and D. Yu, "Self-supervised text-independent speaker verification using prototypical momentum contrastive learning," in *ICASSP*, 2021, pp. 6723–6727.

[13] H. Zhang, Y. Zou, and H. Wang, "Contrastive self-supervised learning for text-independent speaker verification," in *ICASSP*, 2020, pp. 6713–6717.

[14] B. T. Polyak and A. B. Juditsky, "Acceleration of stochastic approximation by averaging," *SIAM journal on control and optimization*, vol. 30, no. 4, pp. 838–855, 1992.

[15] A. Tarvainen and H. Valpola, "Mean teachers are better role models: Weight-averaged consistency targets improve semi-supervised deep learning results," in *NeurIPS*, vol. 30, 2017. [Online]. Available: https://proceedings.neurips.cc/paper/2017/file/68053af2923e00204c3ca7c6a3150cf7-Paper.pdf

[16] J. Villalba, N. Chen, D. Snyder, D. Garcia-Romero, A. McCree, G. Sell, J. Borgstrom, L. P. García-Perera, F. Richardson, R. Dehak, P. A. Torres-Carrasquillo, and N. Dehak, "State-of-the-art speaker recognition with neural network embeddings in nist sre18 and speakers in the wild evaluations," *Comput. Speech Lang.*, vol. 60, 2020.

[17] A. Nagrani, J. S. Chung, W. Xie, and A. Zisserman, "Voxceleb: Large-scale speaker verification in the wild," *Comput. Speech Lang.*, vol. 60, 2020.

[18] J. Thienpondt, B. Desplanques, and K. Demuynck, "The idlab voxceleb speaker recognition challenge 2020 system description," *arXiv:2010.12468*, 2020.

[19] K. He, X. Zhang, S. Ren, and J. Sun, "Deep residual learning for image recognition," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2016.

[20] J. Deng, J. Guo, N. Xue, and S. Zafeiriou, "Arcface: Additive angular margin loss for deep face recognition," in *CVPR*, 2019, pp. 4690–4699.

[21] S. Gao, M.-M. Cheng, K. Zhao, X.-Y. Zhang, M.-H. Yang, and P. H. Torr, "Res2net: A new multi-scale backbone architecture," *IEEE transactions on pattern analysis and machine intelligence*, 2019.

[22] H. Cao, D. G. Cooper, M. K. Keutmann, R. C. Gur, A. Nenkova, and R. Verma, "Crema-d: Crowd-sourced emotional multimodal actors dataset," *IEEE transactions on affective computing*, vol. 5, no. 4, pp. 377–390, 2014.

[23] 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, p. 335, 2008.

[24] R. Lotfian and C. Busso, "Building naturalistic emotionally balanced speech corpus by retrieving emotional speech from existing podcast recordings," *IEEE Transactions on Affective Computing*, 2017.

[25] R. Pappagari, T. Wang, J. Villalba, N. Chen, and N. Dehak, "x-vectors meet emotions: A study on dependencies between emotion and speaker recognition," in *ICASSP*, 2020, pp. 7169–7173.

[26] J. Villalba, N. Chen, D. Snyder, D. Garcia-Romero, A. McCree, G. Sell, J. Borgstrom, F. Richardson, S. Shon, F. Grondin *et al.*, "State-of-the-art speaker recognition for telephone and video speech: The jhu-mit submission for nist sre18." in *Interspeech*, 2019, pp. 1488–1492.