

# The JHU submission to VoxSRC-21: Track 3

Jaejin Cho<sup>1</sup>, Jesús Villalba<sup>1,2</sup>, Najim Dehak<sup>1,2</sup>

<sup>1</sup>Center for Language and Speech Processing, <sup>2</sup>Human Language Technology Center of Excellence,  
Johns Hopkins University, Baltimore, MD, USA

{jcho52, jvillal17, ndehak3}@jhu.edu

## Abstract

This technical report describes Johns Hopkins University speaker recognition system submitted to Voxceleb Speaker Recognition Challenge 2021 Track 3: Self-supervised speaker verification (closed). Our overall training process is similar to the proposed one from the first place team in the last year’s VoxSRC2020 challenge. The main difference is a recently proposed non-contrastive self-supervised method in computer vision (CV), distillation with no labels (DINO), is used to train our initial model, which outperformed the last year’s contrastive learning based on momentum contrast (MoCo). Also, this requires only a few iterations in the iterative clustering stage, where pseudo labels for supervised embedding learning are updated based on the clusters of the embeddings generated from a model that is continually fine-tuned over iterations. In the final stage, Res2Net50 is trained on the final pseudo labels from the iterative clustering stage. This is our best submitted model to the challenge, showing 1.89, 6.50, and 6.89 in EER(%) in voxceleb1\_test\_o, VoxSRC-21 validation, and test trials, respectively.

## 1. Introduction

In Voxceleb Speaker Recognition Challenge 2021 (VoxSRC-21) Track 3: self-supervised speaker verification (SV), a participant is allowed to use only VoxCeleb2 [1] dev subset without any speaker labels for model training. For SV system validation, participants are restricted to use the VoxCeleb1 [2] pairs or the provided pairs composed of the subset of the VoxCeleb1 utterances where their distribution matches that of the test pairs. The validation pairs, however, cannot be used for training. The test pairs contain utterances that are shorter than the utterances in training and validation subsets.

This report shares the details about our developed systems and findings in this challenge. Our main focus was to check how the newly proposed non-contrastive self-supervised learning (SSL) method in CV, DINO [3], works in speaker embedding learning. Thus, the training stages after that follow what’s proposed from the first place team [4] from the last year’s challenge [5] with small modifications. This year’s challenge has a special focus on multi-lingual verification but we do not develop specific systems for this.

## 2. Method

In short, the overall pipeline consists of training a front-end model for speaker embedding extractor and then scoring trial pairs with cosine scoring in the speaker embedding space. In the front-end modeling, we mainly have three stages similar to the first place team’s development pipeline [4] except using non-contrastive learning for the initial model. We explain the three

stages in order.

### 2.1. Initial model training: DINO

To generate pseudo labels from embedding clusters for supervised training, we train an initial model to extract good embeddings that can be clustered by speaker. We adopt a specific non-contrastive SSL method, DINO, for the purpose.

#### 2.1.1. Motivation: Why non-contrastive over contrastive learning

There have been many methods to train embeddings in a self-supervised manner [6, 7, 8, 9, 10, 4, 11, 3], and contrastive loss based methods with data augmentation are popular and well-performing ones [10, 4, 11]. The works using contrastive loss compose the negative samples with different samples from a current sample to make the current and negative samples far from each other in the embedding space. This, however, could be wrong as the size of the queue to accumulate negative samples gets bigger. For example, the more we accumulate utterances to compose negative samples, the more probable some of the utterances are from the same speaker of the current utterance. We could make the queue size smaller to avoid this issue but this degrades the performance [12]. In [11], the author introduced a clustering stage at every epoch to sample negative samples only from different clusters but the improvement was not large.

Non-contrastive methods, however, do not require negative samples so they are free from this issue. Moreover, non-contrastive methods have shown comparable or even better performance compared to contrastive methods. Thus, we propose to apply a non-contrastive SSL method recently proposed, DINO [3], that outperforms previous SSL methods in many CV tasks.

#### 2.1.2. Distillation with No labels

In [3], the author proposed a design to maximize the similarity between feature distributions of differently augmented images from an original image. This is based on the assumption that augmented images from one image keep the same semantic information. For example, although you cropped two images from a dog image and make one a black image while making the other jittered, they are still dog images.

The training is done as follows: First, a given sample is augmented in different ways. To be specific, you crop a local view and a global view of an image, where local and global views mean small and large portions of the image. Several augmentations can be added to the cropped images, such as color jittering, Gaussian blur, solarization, etc. The local views are propagated through one branch while the global views are propagated through the other branch to minimize the cross-entropy between two distributions calculated along the branches. A student net-

Stage	Algorithm/Loss	Model	EER (%)		
			voxceleb1_o_test	VoxSRC-21 val	VoxSRC-21 test
Initial model training (self-supervised learning)	DINO	LResNet34	4.83	13.96	-
	MoCo	ECAPA [4]	7.3	-	-
Iterative clustering	AAM	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) [4]	2.1	-	-
Robust training + larg-margin fine-tuning	AAM	Res2Net50	1.89	6.50	6.88
			1.91	6.32	6.64*

Table 1: Speaker verification results over 3 different trial lists with progressing/different systems over the three stages. The numbers from [4] seems rounded to the nearest tenth. Pseudo labels for robust training were generated from ResNet (iter3) model. \* means the submission happened after the challenge deadline.

work in one branch and a teacher network in the other branch are initialized with the same architecture and the model 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 of the student parameters. To avoid a model to find trivial solutions, i.e., having distributions where one dimension is dominant or uniform distributions, *centering* and *sharpening* are used. *centering* prevents one dimension from dominating by calculating a center by equation. However, using *centering* encourages a uniform distribution. That is why *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. The architecture for the student and teacher networks is composed of a backbone, e.g., ViT [13] or ResNet [14] without later fully connected layers and a projection head. The projection head consists of a 3-layer fully connected layers with hidden dimension 2048 followed by L2 normalization and a weight normalized fully connected layer with  $K$  dimensions.

### 2.1.3. DINO to learn embedding from speech

The assumption that augmented images from one image keep the same semantic information in CV can be similarly applied to speaker embedding learning. For example, most of the speech corpora consist of utterances where each utterance is spoken by one speaker. In this case, it is reasonable to assume that segments extracted from random positions in the same utterance have the same speaker information. The correspondence between CV and speaker embedding learning is following. An image corresponds to an utterance while cropping local and global views from an image corresponds to extracting short and long segments from an utterance. The popular augmentation methods in speaker embedding learning after extracting segments are adding sounds such as babbling, music, noise in the background, or applying room impulse response effects.

## 2.2. Iterative clustering: pseudo label update

In this stage, we train a new model based on pseudo labels generated from the initial model. In detail, we extract speaker embeddings from the initial model and cluster them using clustering algorithms where the number of clusters is heuristically determined based on speaker verification performance on validation data. Indices of final clusters are used as pseudo speaker labels for supervised speaker embedding training.

Once the first labels are generated from the initial model,

we train a new model, possibly with a larger model. The model is continually updated over iterations based on pseudo labels updated after each iteration. The labels are updated in the same way explained above, i.e., through speaker embedding extraction, clustering, and pseudo labeling. The number of clusters is fixed as one value over the iterations.

## 2.3. Robust training on the final pseudo-labels

In this stage, a new larger model is trained with a large margin fine-tuning after a few epochs. The difference from the last year [4] is that we keep using the additive angular margin (AAM) loss instead of sub-center AAM [15] in this stage.

## 3. Experiment and result

For the input features in training, we used an 80-dimensional log filter bank calculated over the 25 ms window with a 10 ms shift. The moving window of 150 ms was used for the mean normalization of the features. Adam [16] optimizer with learning rate scheduling was used over the training stages.

### 3.1. Initial model training: DINO

The embedding architecture we used for the DINO backbone is a light version of ResNet34 (LResNet34) with the kernel size of the first convolution layer as 3 instead of 7 and with a mean and standard deviation pooling layer followed by a fully connected layer to have the embedding dimension as 256 as in [17]. The  $K$  in the following DINO projection head was 65536. The reason for selecting the LResNet34 as the embedding architecture is to reduce the training time considering DINO takes more computation compared to conventional supervised model training. For the augmentation, we added sounds such as babbling, music, noise in the background or applied room impulse response effects.

As shown the Initial model training (self-supervised learning) row in Table 1, DINO outperforms MoCo in the stage.

### 3.2. Iterative clustering: pseudo label update

In this stage, we used an original ResNet34 architecture [14] with the kernel size of the first convolution layer as 3 instead of 7 and with a mean and standard deviation pooling layer followed by a fully connected layer to have the embedding dimension as 256. The loss function used was AAM softmax [18], warming up the margin value from 0 to 0.3 for the first 20 epochs. The number of clusters is set to 7500 following [4].

As shown in the Iterative clustering row in Table 1, the model performance converges from 2nd and 4th iterations on voxceleb1\_test\_o and VoxSRC-21 validation, respectively. This is possibly because the embeddings from DINO initial model are better than the ones from MoCo.

### 3.3. Robust training on the final pseudo-labels

Res2Net50 [19] architecture with 26 for the width of filters and 4 for the scale was used in the final robust training stage. The pooling layer and the following fully connected layers are the same in the previous stages. The AAM loss with the same setting in the previous stage was used. The pseudo labels were generated from the 3rd model from the previous stage, ResNet34 (iter3) in Table 1. After 30 epochs, the post-pooling layers in the model were fine-tuned with a larger margin, 0.5. The large margin fine-tuned model on longer chunks, 4 seconds, however, degraded the performance on VoxSRC-21 test pairs although it showed improvement on validation pairs, as shown in Table 2. This is possibly due to the short-length segments, less than 4-second, that take a large portion of the utterances in the test trials. Thus, we used 3-second chunks instead for large margin fine-tuning.

Segment length	voxceleb1_test_o	VoxSRC-21_val	VoxSRC-21_test
No fine-tuning	<b>1.89</b>	6.50	6.88
2 second	1.97	6.86	-
3 second	1.91	6.32	<b>6.64*</b>
4 second	1.92	<b>6.31</b>	7.23

Table 2: Relationship between segment length and the performance as EER(%) in large margin fine-tuning over 3 different trial lists. \* means the submission happened after the challenge deadline.

Training a larger model and the following large margin fine-tuning improve the performance on VoxSRC-21 validation trials as shown in the Robust training row Table 1.

## 4. Conclusion

We developed speaker verification systems without using speaker labels and achieved 1.91, 6.32, and 6.64 in EER(%) in voxceleb1\_test\_o, VoxSRC-21 validation and test trials, respectively. Our main difference from the previous year’s first place team [4] was to use non-contrastive self-supervised learning method, DINO [3]. This showed better performance in the initial model training stage compared to MoCo [12]. Also, the better speaker embedding in the initial model led to only a few iterations in the next iterative clustering stage. In the robust training stage, we carefully chose the training segment lengths not to overfit to the training or validation subsets while avoiding too short segment lengths. This is because the large portion of the utterances in the test trials was shorter than ones in the training and validation pairs, having the lengths less than 4 seconds.

## 5. References

- [1] J. S. Chung, A. Nagrani, and A. Zisserman, “Voxceleb2: Deep speaker recognition,” in *INTERSPEECH*, 2018.
- [2] A. Nagrani, J. S. Chung, W. Xie, and A. Zisserman, “Voxceleb: Large-scale speaker verification in the wild,” *Comput. Speech Lang.*, vol. 60, 2020.
- [3] M. Caron, H. Touvron, I. Misra, H. Jégou, J. Mairal, P. Bojanowski, and A. Joulin, “Emerging properties in self-supervised vision transformers,” in *Proceedings of the International Conference on Computer Vision (ICCV)*, 2021.
- [4] J. Thienpondt, B. Desplanques, and K. Demuynck, “The idlab voxceleb speaker recognition challenge 2020 system description,” *arXiv preprint arXiv:2010.12468*, 2020.
- [5] A. Nagrani, J. S. Chung, J. Huh, A. Brown, E. Coto, W. Xie, M. McLaren, D. A. Reynolds, and A. Zisserman, “Voxsrc 2020: The second voxceleb speaker recognition challenge,” *arXiv preprint arXiv:2012.06867*, 2020.
- [6] W.-N. Hsu, Y. Zhang, and J. Glass, “Unsupervised learning of disentangled and interpretable representations from sequential data,” in *Proceedings of the 31st International Conference on Neural Information Processing Systems*, 2017, pp. 1876–1887.
- [7] T. Stafylakis, J. Rohdin, O. Plchot, P. Mizera, and L. Burget, “Self-supervised speaker embeddings,” *arXiv preprint arXiv:1904.03486*, 2019.
- [8] Z. Peng, S. Feng, and T. Lee, “Mixture factorized auto-encoder for unsupervised hierarchical deep factorization of speech signal,” in *ICASSP 2020*. IEEE, 2020, pp. 6774–6778.
- [9] J.-B. Grill, F. Strub, F. Althé, 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 *Advances in Neural Information Processing Systems*, H. Larochelle, M. Ranzato, R. Hadsell, M. F. Balcan, and H. Lin, Eds., vol. 33. Curran Associates, Inc., 2020, pp. 21 271–21 284.
- [10] J. Huh, H. S. Heo, J. Kang, S. Watanabe, and J. S. Chung, “Augmentation adversarial training for self-supervised speaker recognition,” *arXiv preprint arXiv:2007.12085*, 2020.
- [11] W. Xia, C. Zhang, C. Weng, M. Yu, and D. Yu, “Self-supervised text-independent speaker verification using prototypical momentum contrastive learning,” in *ICASSP 2021-2021 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 2021, pp. 6723–6727.
- [12] 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.
- [13] A. Dosovitskiy, L. Beyer, A. Kolesnikov, D. Weissenborn, X. Zhai, T. Unterthiner, M. Dehghani, M. Minderer, G. Heigold, S. Gelly, J. Uszkoreit, and N. Houlsby, “An image is worth 16x16 words: Transformers for image recognition at scale,” in *International Conference on Learning Representations*, 2021. [Online]. Available: <https://openreview.net/forum?id=YicbFdNTTy>
- [14] 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.
- [15] J. Deng, J. Guo, T. Liu, M. Gong, and S. Zafeiriou, “Sub-center arcface: Boosting face recognition by large-scale noisy web faces,” in *European Conference on Computer Vision*. Springer, 2020, pp. 741–757.
- [16] D. P. Kingma and J. Ba, “Adam: A method for stochastic optimization,” *arXiv preprint arXiv:1412.6980*, 2014.
- [17] 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.

- [18] J. Deng, J. Guo, N. Xue, and S. Zafeiriou, "Arcface: Additive angular margin loss for deep face recognition," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2019, pp. 4690–4699.
- [19] 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.