



# STC-innovation Speaker Recognition Systems for Far-Field Speaker Verification Challenge 2020

Aleksei Gusev<sup>1,2</sup>, Vladimir Volokhov<sup>2</sup>, Alisa Vinogradova<sup>1,2</sup>, Tseren Andzhukaev<sup>2</sup>,  
Andrey Shulipa<sup>1</sup>, Sergey Novoselov<sup>1,2</sup>, Timur Pekhovsky<sup>2</sup>, Alexander Kozlov<sup>2</sup>

<sup>1</sup>ITMO University, St. Petersburg, Russia

<sup>2</sup>STC-Innovations Ltd., St. Petersburg, Russia

{gusev-a, volokhov, gazizullina, andzhukaev, shulipa,  
novoselov, tim, kozlov-a}@speechpro.com

## Abstract

This paper presents speaker recognition (SR) systems submitted by the Speech Technology Center (STC) team to the Far-Field Speaker Verification Challenge 2020. SR tasks of the challenge are focused on the problem of far-field text-dependent speaker verification from single microphone array (Track 1), far-field text-independent speaker verification from single microphone array (Track 2) and far-field text-dependent speaker verification from distributed microphone arrays (Track 3).

In this paper, we present techniques and ideas underlying our best performing models. A number of experiments on x-vector-based and ResNet-like architectures show that ResNet-based networks outperform x-vector-based systems. Submitted systems are the fusions of ResNet34-based extractors, trained on 80 Log Mel-filter bank energies (MFBs) post-processed with U-net-like voice activity detector (VAD). The best systems for the Track 1, Track 2 and Track 3 achieved 5.08% EER and 0.500  $C_{det}^{min}$ , 5.39% EER and 0.541  $C_{det}^{min}$  and 5.53% EER and 0.458  $C_{det}^{min}$  on the challenge evaluation sets respectively.

**Index Terms:** FFSVC, speaker recognition, deep neural network, domain adaptation, neural network-based VAD.

## 1. Introduction

Text-independent speaker verification is one of the most significant directions in speech processing and biometric authentication areas. The growing market of voice-controlled smart devices is driving the demand for reliable far-field speech and speaker recognition technology. The Voices Obscured in Complex Environmental Settings (VOiCES) from a Distance Challenge 2019 [1, 2, 3] was held to promote research in the area of speaker recognition and automatic speech recognition (ASR) for far-field speech recorded under noisy conditions. The Short-duration Speaker Verification (SdSV) Challenge 2020 [4] gives more evidence of increasingly serious scientific interest towards text-independent speaker verification techniques. It sheds the light on another property of speaker verification systems, that promotes speaker verification research and influences the growth of the voice assistant market – the ability to act reliably in a short duration scenario.

Far-Field Speaker Verification Challenge 2020 (FFSVC 2020) [5, 6] heads in the same direction with the special focus on far-field distributed microphone arrays under noisy conditions in real scenarios. The most challenging part of the contest setup is cross-channel testing.

The main focus of this paper is the Track 2 problem of the FFSVC 2020 Challenge. We have also applied several hypotheses tested on this track to the other two tracks as well. Following

the success of deep neural network (DNN)-based speaker embedding extractors [7, 8, 9], we built our systems on the state-of-the-art DNN topologies.

A TDNN-based x-vector system significantly outperforms conventional i-vector-based systems in terms of speaker recognition performance and hence became a new baseline for text-independent SR tasks [7]. Thus, the first topology we have tried for the frame-level part of our system is extended Time Delay Neural Network (TDNN) [10], which processes global spectral and local temporal information and gradually expands context over time.

Earlier works on short utterance speaker recognition and far-field speaker verification [8, 9] demonstrate that substantial improvements can be achieved by deeper architectures such as residual networks [11]. Inspired by these papers we also used ResNet for the frame-level part of the embedding extractor network. Different from TDNN-based approaches, a ResNet-like network is deeper and applies 2-dimensional convolutions over the input features. This way both local temporal and local spectral information is attended to with an equal priority.

The main results on the challenge tasks are obtained by ResNet34-based extractors. In the context of studied typologies, the extraction of features simultaneously from spectral and time dimensions leads to significantly better results compared with the work of 1-dimensional convolutions over the time dimension. Also, our experiments show that additional improvements of recognition performance can be obtained by fine-tuning embedding extractors with the in-domain Chinese Mandarin datasets, preprocessing data with more accurate voice activity detection at the training stage. Gains in speaker recognition quality can also be achieved by the use of various adaptation techniques at the training and testing stages, as well as by the utilization of multi-channel fusion approaches at the testing stage.

The paper is structured as follows: Section 2 describes the feature extraction procedure, voice activity detection techniques used to preprocess data and architectures used for embedding extraction. It also outlines domain adaptation, channel fusion and backend strategies. Section 3 introduces the datasets used to train and test the presented systems. Experimental design and results are presented in Sections 4 and 5 respectively.

## 2. System components

### 2.1. Feature extraction

Input features for all systems presented in this paper are 80-dimensional Log Mel-filter bank energies extracted from 16kHz raw input signals. We compute MFBs from the signal with

25ms frame-length and 15ms overlap.

Additionally, we use per-utterance Cepstral Mean Normalization (CMN) over a 3-second sliding window over the stack of MFBs to compensate for the channel effects and noise by transforming data to have zero mean [12]. The VAD was used after the CMN-normalization procedure. We further apply global mean and standard deviation (std) normalization for each utterance with the pre-computed 80-dimensional vectors of means and stds over this utterance.

## 2.2. Voice activity detector

In this work we explored two types of VADs for the SR task:

- energy-based VAD from the Kaldi Toolkit [13];
- neural network-based VAD.

The Neural network-based VAD uses U-net architecture [14] as a backbone and is described in more detail in our papers [15, 16]. It was trained on the large-scale dataset of telephone channel audios (telephone part of data from the NIST SRE challenge and our proprietary Russian speech subcorpus RusTelecom [17]). We have found that the VAD trained on telephone data produces high-quality results for the microphone channel audios, that is the reason why we did not train our VAD for microphone data from scratch.

Speech labels for the VAD training were obtained by manual segmentation and using ASR-based VAD processing of a clean version of the data. Preprocessing of VAD input data consisted of extraction of 8kHz 23-dimensional Mel-frequency cepstral coefficient (MFCC) features from the raw signal with 25ms frame-length and 20ms shift. We found 23-dimensional MFCCs to be a tradeoff between the quality of embeddings extracted from these features and the speed of training and inference.

Since training and test data used for the FFSVC 2020 consists of 16kHz microphone speech, VAD markup was first extracted from the 23-dimensional MFCC features, computed for the audios down-sampled from 16kHz to 8kHz, then the resultant markup was used to extract voiced frames from the 80-dimensional MFB features calculated from the same raw waveform data.

## 2.3. Embedding extractors

We used two kinds of neural network architectures to process acoustic features – ResNet-like and x-vector-based systems. The main conceptual difference between them is the way convolutions operate on time and spectral dimensions. ResNets utilize both local spectral and time information to learn features for the output feature maps, while x-vector-based systems use 1-dimensional convolutions over the time domain and thus utilize global spectral and local time information.

**ResNet-like.** Table 1 describes the ResNet34 [11] architecture we used. ReLU activation and batch normalization follow each convolutional layer, and Maxout activation [18] is used for the embedding layer. The statistics pooling layer aggregates the input features of the segment-level embedding block over time and spectral dimensions. In the paper, we used several embedding extractors based on ResNet34 architecture. These extractors differ in the data used for the training, the type of data augmentation applied, and the voice activity detector used to preprocess input.

**X-vector-based.** We took extended TDNN-based x-vectors [10] as a baseline. Then we removed dilations from convolutional kernels to avoid the possibility for the gridding artifacts

Table 1: Architecture configuration of the embedding extractor based on ResNet34

Layer name	Structure	Output
Input	80 MFB log-energy	$80 \times 200 \times 1$
Conv2D-1	$3 \times 3$ , stride 1	$80 \times 200 \times 32$
ResNet-1	$\begin{bmatrix} 3 \times 3, 32 \\ 3 \times 3, 32 \end{bmatrix} \times 3$ , st. 1	$80 \times 200 \times 32$
ResNet-2	$\begin{bmatrix} 3 \times 3, 64 \\ 3 \times 3, 64 \end{bmatrix} \times 4$ , st. 2	$40 \times 100 \times 64$
ResNet-3	$\begin{bmatrix} 3 \times 3, 128 \\ 3 \times 3, 128 \end{bmatrix} \times 6$ , st. 2	$20 \times 50 \times 128$
ResNet-4	$\begin{bmatrix} 3 \times 3, 256 \\ 3 \times 3, 256 \end{bmatrix} \times 3$ , st. 2	$10 \times 25 \times 256$
StatsPool	mean and std	$20 \times 256$
Flatten	–	5120
Dense1	embedding layer	512
Dense2	output layer	$N_{spk}$

and replaced ReLU activation with its leaky version (Table 2). Several experiments with the width and depth of the network showed no further improvement for the number of filters beyond 512 and the number of frame-level blocks (comprised of 1-dimensional convolutional layer and  $1 \times 1$  1-dimensional convolutional layer) beyond 4. Similarly, the addition of extra fully connected segment-level layers did not affect the verification accuracy. It was found that after the network reaches a certain depth, the performance tends to saturate towards mean accuracy due to the lack of local spectral and global temporal information in the intermediate layers. We tried to add temporal attention based on the Squeeze-and-Excitation (SE) block [19] after each convolutional block to overcome the latter limitation, however, no significant improvement over the basic version was reached. Further experiments with the 2D convolutional networks (ResNets) demonstrate the need for the solution of former limitation together with the transformations of the time axis.

Table 2: Architecture configuration of the embedding extractor based on x-vectors

Layer name	Layer context	Output
Input	80 MFB log-energy	$80 \times 200$
Frame1.1	$[t - 2 : t - 2]$	$512 \times 200$
Frame1.2	$[t]$	$512 \times 200$
Frame2.1	$[t - 1 : t - 1]$	$512 \times 200$
Frame2.2	$[t]$	$512 \times 200$
Frame3.1	$[t - 1 : t - 1]$	$512 \times 200$
Frame3.2	$[t]$	$512 \times 200$
Frame4.1	$[t - 1 : t - 1]$	$512 \times 200$
Frame4.2	$[t]$	$512 \times 200$
Frame5	$[t]$	$1500 \times 200$
StatsPool	mean and std	$2 \times 1500$
Flatten	–	3000
Dense1	embedding layer	512
Dense2	output layer	$N_{spk}$

## 2.4. Domain adaptation

In this work, we utilized different domain adaptation techniques:

- based on the addition of the in-domain data to the train-

ing set to fine-tune and train embedding extractor that solves close-set speaker identification task;

- based on mean speaker embedding subtraction. The mean vector is calculated over the training FFSVC dataset;
- based on two mean speaker embedding subtraction. The main idea is to calculate two vectors of mean values over the training FFSVC dataset for the enrollment and test files independently;
- based on the MultiReader adaptation technique [20], which was used at the embedding extractor training stage. The main idea is to train embedding extractor using two heads, each of which is intended to classify either a large number of speaker IDs from out-of-domain data or a small number of speaker IDs from in-domain data.

### 2.5. Multi-channel fusion

Trial pairs are constructed of a single enrollment recording from a 25cm distance cell phone and multiple test recordings from a single (Track 1 and Track 2) or multiple (Track 3) far-field microphone arrays. The presence of multiple test files for the enrolled speaker fragment allowed us to fuse information from test utterances in several ways:

- by averaging all enrollment-test trials scores for each trial;
- by choosing the maximum score from comparison set of one enrollment embedding with all test embeddings for each trial;
- by computing average test embedding for one trial and comparing it with the associated enrollment embedding.

We found that embedding averaging works slightly better than other methods in terms of verification metrics.

### 2.6. Back-end scoring

Cosine similarity was chosen to be used as a back-end scoring method:

$$S(\mathbf{x}_1, \mathbf{x}_2) = \frac{\mathbf{x}_1^T \mathbf{x}_2}{\|\mathbf{x}_1\| \|\mathbf{x}_2\|}, \quad (1)$$

where  $(\mathbf{x}_1, \mathbf{x}_2)$  are speaker embedding vectors.

## 3. Datasets

We used several datasets for our experiments.

**Model pre-training dataset.** We used concatenated VoxCeleb1 and VoxCeleb2 (SLR47) [21] corpus datasets to pre-train all our models. The overall number of speakers in the resultant set was 7146. Augmented data was generated using standard Kaldi augmentation recipe (reverberation, babble, music and noise) using the freely available MUSAN and simulated Room Impulse Response (RIR) datasets<sup>1</sup>.

**Model fine-tuning dataset.** We have expanded model pre-training data with additional 2099 speakers from several Chinese Mandarin corpora (SLR33, SLR62, SLR82, SLR85, FFSVC train set) to add domain knowledge. We concatenated multiple short-duration utterances of one speaker into multiple larger files of 20 sec for training convenience. SLR33 and SLR62 sets were augmented similarly to SLR47 set, whereas

SLR82 and SLR85 relatively noisy sets were augmented only with reverberation. Both augmented and non-augmented versions of the FFSVC train set were used in our experiments. During the construction of the extended dataset speaker overlaps between different datasets were taken into account.

**Development and evaluation data.** FFSVC 2020 challenge development set consisting of 35 speakers was used for verification systems testing. Recordings from five devices (one close-talk microphone, one 25cm distance cellphone and three randomly selected microphone arrays) for each utterance are provided. The submission system uses the FFSVC evaluation set for the final testing. It includes 80 speakers and enrollment and test recordings are done by different types of devices. Detailed descriptions of the challenge data are given in [6].

## 4. Experiments

### 4.1. Pre-training of embedding extractor

Both ResNet and x-vector-based embedding extractor were trained on model pre-training dataset from Section 3 We performed training of our models using batches, consisting of randomly sampled sequences of 200 MFB features (2s of speech). AM-Softmax loss function [22] was taken for the objective (with optimal margin and scale parameter settings fixed to 0.2 and 30 respectively). For ResNet-like models we used Adam optimizer with the starting learning rate fixed to 0.001 and kept dividing it by ten every second epoch. In x-vector training, SGD demonstrated better convergence in the combination with cyclic learning rate scheduling policy with the minimum and maximum learning rate parameters set to 0.002 and 0.12 respectively. During one epoch the full pass of SLR47 data was done.

### 4.2. Fine-tuning of embedding extractors

Fine-tuning of the ResNet-like extractors was done in two steps. First, we trained segment-level layers and newly initialized classification head, with convolutional layers frozen. Second, we unfroze convolutional layers and re-trained the overall network with a low learning rate. We used the FFSVC development set and original VoxCeleb 1 test set for model validation in this task. Adam optimizer was used for model optimization, AM-Softmax loss function with margin and scale parameter settings fixed to 0.2 and 30 respectively was taken.

The MultiReader adaptation technique [20] was used as an alternative fine-tuning approach. We trained an embedding extractor based on ResNet34 architecture using two heads. The first head is used to classify a large number of speaker IDs from out-of-domain data (7146 speakers from VoxCeleb1 and VoxCeleb2 datasets). The second head is used to classify a small number of speaker IDs from the in-domain data (120 speakers from the FFSVC train set). Since the in-domain dataset does not contain enough amount of data, training the embedding extractor on it can lead to overfitting. Requiring the embedding extractor to perform reasonably well also on out-of-domain data helps to regularize the embedding extractor. We tried to minimize the following cost function based on AM-Softmax:

$$\mathcal{L}(D; \mathbf{W}) = 0.5 \cdot \mathcal{L}(D_1; \mathbf{W}_1) + 0.5 \cdot \mathcal{L}(D_2; \mathbf{W}_2), \quad (2)$$

where  $D_1$  and  $D_2$  are out-of-domain and in-domain data correspondingly,  $\mathbf{W}_1$  and  $\mathbf{W}_2$  are model parameter sets that allow to compute outputs for the first and second heads correspondingly,  $D = D_1 \cup D_2$ ,  $\mathbf{W} = \mathbf{W}_1 \cup \mathbf{W}_2$ .

For the initialization of the frame-level, segment-level and the first output layers of the two-headed ResNet34 model we

<sup>1</sup><http://www.openslr.org>

Table 3: Results of our systems on FFSVC 2020 Track 2 (development set)

ID	System	Properties	$C_{det}^{min}$	EER
1	ResNet34	initial model, energy VAD, max score for test files	0.820	8.95
2	ResNet34	ID1 + neural VAD	0.712	8.23
3	ResNet34	ID2 + mean vector for test files	0.700	8.24
4	ResNet34	ID3 + more data after VAD	0.688	8.26
5	ResNet34	ID4 + mean vector subtraction	0.672	8.22
6	ResNet34	ID5 + two mean vector subtraction	0.655	7.47
7	ResNet34	ID4 + fine-tune model	0.562	4.85
8	ResNet34	ID7 + two mean vector subtraction	0.484	4.46
9	X-vectors	Extended TDNN + energy VAD	0.890	12.30
10	ResNet34	ID1 + MultiReader, mean vector for test files	0.627	5.46

used the single-headed ResNet34 model [15] trained on original and augmented VoxCeleb1 and VoxCeleb2 datasets. The first layer and the frame level of the two-headed ResNet34 model were frozen at the beginning of the MultiReader adaptation procedure. Layer freezing was maintained until convergence and then all layers were unfrozen and training procedure was continued with a reduced learning rate until convergence. We used the original and Kaldi augmented FFSVC train set in this case.

## 5. Results and discussion

Table 3 displays results of the experiments on several systems developed for the Track 2 of the challenge. The performance is measured on the FFSVC development set in terms of EER (Equal Error Rate) and  $C_{det}^{min}$  (Minimum Detection Cost).

The ResNet34-based system with energy VAD and no adaptation to the specificity of the domain was taken for the baseline. We used only concatenated VoxCeleb1 and VoxCeleb2 datasets and its augmented versions for training our baseline system. The maximum score fusion method mentioned in Subsection 2.5 was used to average information from several multi-microphone test utterances. A number of changes were done to the baseline system:

- the expansion of the training set and speaker IDs with Chinese Mandarin datasets;
- the use of U-net-like VAD in training and testing stages;
- the use of one of the adaptation methods described in Subsection 2.4;
- computation of average test embedding vector for one trial and comparison of it with the associated enrollment embedding.

Analysis of results allows us to make the following conclusions based on Table 3:

Table 4: Results of our best single and fused systems on FFSVC 2020 (development set and evaluation set)

System	Task	Dev set		Eval set	
		$C_{det}^{min}$	EER	$C_{det}^{min}$	EER
Single-1	Track 1	0.490	4.24	–	–
Fusion-1	Track 1	0.462	3.66	0.500	5.08
Single-2	Track 2	0.484	4.46	0.564	5.61
Fusion-2	Track 2	0.472	4.28	0.541	5.39
Single-3	Track 3	0.434	3.35	–	–
Fusion-3	Track 3	0.417	3.22	0.458	5.53

- expansion of the training set with the in-domain Chinese Mandarin datasets allows to improve performance of verification system (ID7);
- the use of U-net-like VAD at training and testing stages allows to improve the performance of verification system. This gain is attributable to the fact that U-net-like VAD produces more accurate speech detections in the presence of distortions compared to energy-based VAD (ID2);
- adaptation by one or two mean speaker embeddings subtraction allows to improve the performance of the verification system. The two mean adaptation technique gives special improvement as it accounts for the fact that enrollment and test recordings are formed using various devices (ID6, ID8);
- the use of the MultiReader adaptation technique improves the performance of the verification system. However, this approach requires careful selection of the learning rate (ID10);
- the use of multi-channel fusion approaches gives better results than the comparison between enrollment embedding and randomly selected test embedding for one trial. We did not see much difference in the performance of the verification system for different multi-channel fusion approaches (ID2, ID3).

We presented the best results of our single and fused verification systems for all tracks in the Table 4. We used approaches similar to Track 2 to improve the performance of the verification system for Track 1 and Track 3. However, we used only the text-dependent part of the FFSVC train set for adaptation in Track 1 and Track 3 for better performance of our systems.

## 6. Conclusions

Obtained results confirm that deep ResNet architectures are robust and allow to obtain a good quality of speaker verification for short-duration utterances. Our best performing system for FFSVC 2020 (development set) protocols is ResNet34-based system built on high-frequency resolution MFB features. It is trained with AM-Softmax-based loss function. We should also note that utilization of additional in-domain data, our U-net-like VAD, various adaptation techniques, and multi-channel fusion approaches provide additional performance gains for proposed SR systems in considered tasks.

## 7. Acknowledgements

This work was partially financially supported by the Government of the Russian Federation (Grant 08-08) and by the Foundation NTI (contract 20/18gr) ID 000000007418QR20002.

## 8. References

- [1] M. K. Nandwana, J. van Hout, C. Richey, M. McLaren, M. A. Barrios, and A. Lawson, "The VOICES from a distance challenge 2019," in *INTERSPEECH 2019 – 20<sup>th</sup> Annual Conference of the International Speech Communication Association, September 15-19, Graz, Austria, Proceedings*, 2019, pp. 2438–2442.
- [2] C. Richey, M. A. Barrios, Z. Armstrong, C. Bartels, H. Franco, M. Graciarena, A. Lawson, M. K. Nandwana, A. Stauffer, J. van Hout, P. Gamble, J. Hetherly, C. Stephenson, and K. Ni, "Voices obscured in complex environmental settings (VOICES) corpus," in *INTERSPEECH 2018 – 19<sup>th</sup> Annual Conference of the International Speech Communication Association, September 2-6, Hyderabad, India, Proceedings*, 2018, pp. 1566–1570.
- [3] S. Novoselov, A. Gusev, A. Ivanov, T. Pekhovsky, A. Shulipa, G. Lavrentyeva, V. Volokhov, and A. Kozlov, "STC speaker recognition systems for the VOICES from a distance challenge," in *INTERSPEECH 2019 – 20<sup>th</sup> Annual Conference of the International Speech Communication Association, September 15-19, Graz, Austria, Proceedings*, 2019, pp. 2443–2447.
- [4] H. Zeinali, K. A. Lee, J. Alam, and L. Burget, "Short-duration speaker verification (SdSV) challenge 2020: the challenge evaluation plan," *arXiv preprint arXiv:1912.06311*, 2019.
- [5] X. Qin, M. Li, H. Bu, R. K. Das, W. Rao, S. Narayanan, and H. Li, "The FFSVC 2020 evaluation plan," *arXiv preprint arXiv:2002.00387*, 2020.
- [6] X. Qin, M. Li, H. Bu, W. Rao, R. K. Das, S. Narayanan, and H. Li, "The INTERSPEECH 2020 far-field speaker verification challenge," in *to appear in INTERSPEECH 2020 – 21<sup>th</sup> Annual Conference of the International Speech Communication Association, October 25-29, Shanghai, China, Proceedings*, 2020.
- [7] D. Snyder, D. Garcia-Romero, D. Povey, and S. Khudanpur, "Deep neural network embeddings for text-independent speaker verification," in *INTERSPEECH 2017 – 18<sup>th</sup> Annual Conference of the International Speech Communication Association, August 20-24, Stockholm, Sweden, Proceedings*, 2017, pp. 999–1003.
- [8] A. Hajavi and A. Etemad, "A deep neural network for short-segment speaker recognition," in *INTERSPEECH 2019 – 20<sup>th</sup> Annual Conference of the International Speech Communication Association, September 15-19, Graz, Austria, Proceedings*, 2019, pp. 2878–2882.
- [9] Q. Xiaoyi, D. Cai, and M. Li, "Far-field end-to-end text-dependent speaker verification based on mixed training data with transfer learning and enrollment data augmentation," in *INTERSPEECH 2019 – 20<sup>th</sup> Annual Conference of the International Speech Communication Association, September 15-19, Graz, Austria, Proceedings*, 2019, pp. 4045–4049.
- [10] D. Garcia-Romero, D. Snyder, G. Sell, A. McCree, D. Povey, and S. Khudanpur, "X-vector DNN refinement with full-length recordings for speaker recognition," in *INTERSPEECH 2019 – 20<sup>th</sup> Annual Conference of the International Speech Communication Association, September 15-19, Graz, Austria, Proceedings*, 2019, pp. 1493–1496.
- [11] K. He, X. Zhang, S. Ren, and J. Sun, "Deep residual learning for image recognition," in *CVPR 2016 – 2016 IEEE Conference on Computer Vision and Pattern Recognition, June 26 - July 1, Las Vegas, Nevada, USA, Proceedings*, 2016, pp. 770–778.
- [12] S. Furui, "Cepstral analysis technique for automatic speaker verification," *IEEE Transactions on Acoustics, Speech, and Signal Processing*, vol. 29, no. 2, pp. 254–272, 1981.
- [13] D. Snyder, D. Garcia-Romero, G. Sell, D. Povey, and S. Khudanpur, "X-vectors: Robust DNN embeddings for speaker recognition," in *ICASSP 2018 – IEEE International Conference on Acoustics, Speech and Signal Processing, April 15-20, Calgary, Canada, Proceedings*, 2018, pp. 5329–5333.
- [14] O. Ronneberger, P. Fischer, and T. Brox, "U-net: Convolutional networks for biomedical image segmentation," in *MICCAI 2015 – 18<sup>th</sup> International Conference on Medical Image Computing and Computer Assisted Intervention, October 5-9, Munich, Germany, Proceedings*, 2015, pp. 234–241.
- [15] A. Gusev, V. Volokhov, T. Andzhukaev, S. Novoselov, G. Lavrentyeva, M. Volkova, A. Gazizullina, A. Shulipa, A. Gorlanov, A. Avdeeva, A. Ivanov, A. Kozlov, T. Pekhovsky, and Y. Matveev, "Deep speaker embeddings for far-field speaker recognition on short utterances," in *Odyssey 2020 – The Speaker and Language Recognition Workshop, November 02-05, Tokyo, Japan, Proceedings*, 2020.
- [16] G. Lavrentyeva, M. Volkova, A. Avdeeva, S. Novoselov, A. Gorlanov, T. Andzhukaev, A. Ivanov, and A. Kozlov, "Blind speech signal quality estimation for speaker verification systems," in *to appear in INTERSPEECH 2020 – 21<sup>th</sup> Annual Conference of the International Speech Communication Association, October 25-29, Shanghai, China, Proceedings*, 2020.
- [17] S. Novoselov, A. Shulipa, I. Kremnev, A. Kozlov, and V. Shchemelinin, "On deep speaker embeddings for text-independent speaker recognition," in *Odyssey 2018 – The Speaker and Language Recognition Workshop, June 26-29, Les Sables d'Olonne, France, Proceedings*, 2018, pp. 378–385.
- [18] I. J. Goodfellow, D. Warde-Farley, M. Mirza, A. Courville, and Y. Bengio, "Maxout networks," in *ICML 2013 – 30<sup>th</sup> International Conference on Machine Learning, June 17-19, Atlanta, Georgia, USA, Proceedings*, vol. 28, no. 3, 2013, pp. 1319–1327.
- [19] J. Hu, L. Shen, and G. Sun, "Squeeze-and-excitation networks," in *CVPR 2018 – 2018 IEEE Conference on Computer Vision and Pattern Recognition, June 18-23, Salt Lake City, Utah, USA, Proceedings*, 2018, pp. 7132–7141.
- [20] L. Wan, Q. Wang, A. Papir, and I. Moreno, "Generalized end-to-end loss for speaker verification," in *ICASSP 2018 – IEEE International Conference on Acoustics, Speech and Signal Processing, April 15-20, Calgary, Canada, Proceedings*, 2018, pp. 4879–4883.
- [21] A. Nagrani, J. S. Chung, W. Xie, and A. Zisserman, "VoxCeleb: Large-scale speaker verification in the wild," *Computer Speech and Language*, vol. 60, p. 101027, 2020.
- [22] F. Wang, W. Liu, H. Liu, and J. Cheng, "Additive margin softmax for face verification," *IEEE Signal Processing Letters*, vol. 25, no. 7, pp. 926–930, 2018.