

LEARNING DOMAIN INVARIANT REPRESENTATIONS FOR CHILD-ADULT CLASSIFICATION FROM SPEECH

Rimita Lahiri¹, Manoj Kumar¹, Somer Bishop², Shrikanth Narayanan¹

¹Signal Analysis and Interpretation Laboratory, University of Southern California

²Department of Psychiatry, University of California, San Francisco

ABSTRACT

Diagnostic procedures for ASD (*autism spectrum disorder*) involve semi-naturalistic interactions between the child and a clinician. Computational methods to analyze these sessions require an end-to-end speech and language processing pipeline that go from raw audio to clinically-meaningful behavioral features. An important component of this pipeline is the ability to automatically detect who is speaking when i.e., perform child-adult speaker classification. This binary classification task is often confounded due to variability associated with the participants' speech and background conditions. Further, scarcity of training data often restricts direct application of conventional deep learning methods. In this work, we address two major sources of variability—age of the child and data source collection location—using domain adversarial learning which does not require labeled target domain data. We use two methods, generative adversarial training with inverted label loss and gradient reversal layer to learn speaker embeddings invariant to the above sources of variability, and analyze different conditions under which the proposed techniques improve over conventional learning methods. Using a large corpus of *ADOS-2* (*autism diagnostic observation schedule, 2nd edition*) sessions, we demonstrate upto 13.45% and 6.44% relative improvements over conventional learning methods.

Index Terms— Child speech, domain adversarial learning, gradient reversal, autism spectrum disorder

1. INTRODUCTION

Autism spectrum disorder (ASD) refers to a group of neurodevelopmental disorders characterized by abnormalities in speech and language [1, 2, 3] and often diagnosed in children using semi-structured dyadic interactions with a trained clinician. The reported ASD prevalence has been steadily increasing among children in the US: from 1 in 150 [4] to 1 in 59 [5], Computational processing of the participants' speech and language during such child-adult interactions has shown potential in recent years in supporting and augmenting human perceptual and decision making capabilities. [6, 7, 8].

However, previous works utilized manual speaker labels and transcripts for behavioral feature computation, which can be expensive and time-consuming to create. Hence, feature extraction at-scale is dependent on a robust speech and lan-

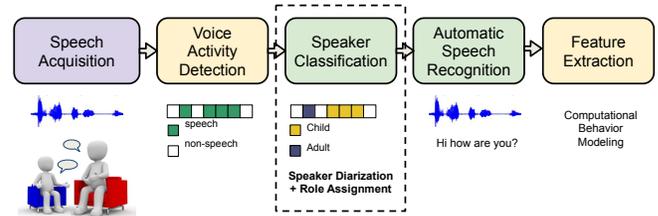


Fig. 1: Speech processing pipeline for feature extraction

guage pipeline (Figure 1). An important component of the pipeline is speaker diarization, which answers the question “*who spoke when?*”. In the context of ASD diagnostic assessment sessions, diarization can be approached as (supervised) child-adult speaker classification. Training a child-adult classification system is often not straightforward due to multiple sources of variability in the data. Among others, two primary sources of variability arise from developmental aspects of child speech [9] and from varying background conditions, often influenced by where and how the data are collected. In this work, we train a child-adult classification system using domain adversarial training [10, 11] to address these sources of variability.

A generative adversarial network (GAN) is composed of two mutually pitting neural networks, termed as the generator and the discriminator. These networks play a minimax game, where the generator aims to create fake samples from a noise vector of some arbitrary distribution in order to confuse the discriminator. On the other hand, the discriminator tries to distinguish between the real and fake samples. Domain adversarial learning can be formulated as a variant of GANs, where the noise vectors are replaced with target data, and the (domain) discriminator network tries to discriminate whether a sample belongs to source or target domain. Hence, the generator network learns to extract domain-invariant representations. The speaker classifier is trained on the generator outputs in a multi-task manner. In this work, we have used two different methods of domain adversarial training namely *Gradient Reversal (GR)* [11] and *Generative Adversarial Networks (GAN)* [12]. GR tries to learn the domain-invariant feature by reversing the gradients coming from domain discriminator while GAN aims to achieve the same by training with inverted domain labels. The full network configuration comprising of generator (feature extractor), discriminator and speaker classifier is shown in Figure 2.

The rest of the paper is organized as follows: Section 2 provides a brief overview of the background works. Section 3 describes the domain adversarial methods used in this work. Section 4 provides experimental details and details of the dataset used. Key outcomes of the experiments are tabulated and interpreted in section 5. Finally, section 6 provides conclusions and highlights possible future extensions.

2. BACKGROUND

2.1. Speaker Diarization in Autism Diagnosis Sessions

Although there exists a significant amount of work in speaker diarization of broadcast news and meetings, interest in spontaneous and real-life conversations has emerged only recently. Diarization solutions for child speech (both child-directed and adult-directed) initially looked at traditional feature representations (MFCCs, PLPs) [13] and speaker segmentation/clustering methods (generalized likelihood ratio, Bayesian information criterion) [14, 15]. In [14], the authors introduced several methods for working with audio collected from children with autism using a wearable device. More recently, approaches based on fixed-dimensional embeddings such as ivectors [16] and DNN speaker embeddings such as x-vectors [17] were explored. While some of the above approaches have adapted clustering methods to child speech [17], to the best of our knowledge none of them have taken into account shifts in domain distribution that is likely to adversely impact diarization performance.

2.2. Domain Adversarial Learning

Domain adaptation within adversarial learning was first introduced by [11] for computer vision related applications. Since then, there has been an emerging trend to use domain adversarial learning to alleviate the mismatch between the training and testing data in various speech applications including ASR and acoustic emotion recognition [18]. In [19, 20] the authors have employed domain adversarial training to improve the robustness of the speech recognition system to handle different noise types and levels. In [21], the authors applied domain adversarial training to address mismatch between close-talk and single-channel far-field recordings. Our motivation for applying domain adversarial learning is inspired from recent applications ([12, 22]) in speaker verification across multiple languages. It was shown that adversarial training can be used to learn robust speaker embeddings across different conditions. We extend this concept to the task of child-adult classification from speech, where variabilities in children’s linguistic capabilities and recording locations can be viewed as domain shift that can be modeled using adversarial learning.

3. DOMAIN ADVERSARIAL LEARNING FOR SPEAKER CLASSIFICATION

The main aim of the work is to efficiently distinguish between the speakers (namely, child and an adult interlocutor) from

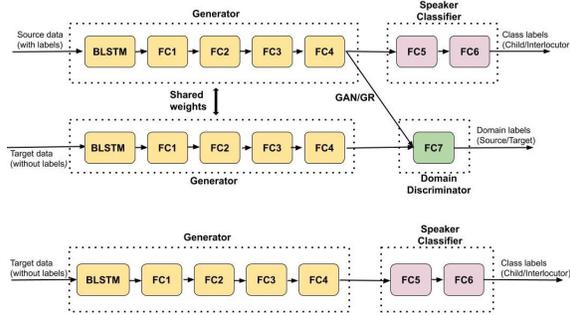


Fig. 2: Training and Testing Network Architecture

audio recordings of diagnostic sessions from different clinical locations. Besides learning domain invariant features by confusing the discriminator, the network must be able to efficiently distinguish between the speakers as well. In this work, we have shown that the proposed objective can be accomplished using a GAN based method, or a GR based method.

Consider samples from the source domain $(X_s, Y_s) \in \Omega_s$ and target domain $(X_t, Y_t) \in \Omega_t$ with a common label space Y . During training, labels from the target domain are assumed unavailable, and data distributions of X_s and X_t might differ. The goal of domain adversarial learning is to maximize the target accuracy by jointly learning to maximize task performance and reducing domain shift between the source and target domains in generator output embedding space.

In our work, we begin by training the network with source data and corresponding speaker labels to minimize task loss. We refer to this as *pre-training*. Following, the adversarial game continues where the discriminator is trained with true domain labels and the generator is trained either with inverted domain labels (*GAN*) or reverse gradients (*GR*) alternatively until convergence is reached.

In both methods, for every batch of data, the training is carried out in three distinct steps. In the first step, the generator and speaker classification models are trained with true speaker labels from the source data using the following objective:

$$\min_{G,C} Loss_{Spk}(X_s, Y_s) = \mathbb{E}_{x_s, y_s \sim (X_s, Y_s)} \sum_{k=1}^2 \mathbb{1}_{k=y_s} \log(C(G^s(x_s))) \quad (1)$$

where $G(\cdot)$ and $C(\cdot)$ are the generator and classifier functions, respectively. In the second step, the embeddings are extracted from the output layer of the generator for both source and target data using the model trained in the previous step. The domain discriminator is now trained with the true domain labels. This step ensures that the discriminator is well trained to distinguish between source and target domain.

$$\min_D Loss_{Dom}(X_s, X_t, G) = \mathbb{E}_{x_s \sim X_s} \log(D(G(x_s))) + \mathbb{E}_{x_t \sim X_t} \log(1 - (D(G(x_t)))) \quad (2)$$

Table 1: Demographic details of ADOS dataset

Category	Statistics
Age(years)	Range: 3.58-13.17 (mean,std):(8.61,2.49)
Gender	123 male, 42 female
Non-verbal IQ	Range: 47-141 (mean,std):(96.01,18.79)
	86 ASD,42 ADHD
Clinical Diagnosis	14 mood/anxiety disorder 12 language disorder 10 intellectual disability, 1 no diagnosis
Age distribution	Cincinnati: ≤ 5 yrs 7, 5-10 yrs 52, ≥ 10 yrs 25 Michigan: ≤ 5 yrs 11, 5-10 yrs 42, ≥ 10 yrs 28

The first and second steps are the same for both GAN and GR: they differ in the third step. For GAN, the generator is trained with source and target data but with inverted domain labels:

$$\min_G \text{Loss}_{Adv}(X_s, X_t, G) = \mathbb{E}_{x_s \sim X_s} \log(D(G(x_t))) + \mathbb{E}_{x_t \sim X_t} \log(1 - (D(G(x_s)))) \quad (3)$$

In case of GR, the gradients from the domain discriminator are reversed for training. In both the cases, the final step ensures the generator is trained well to generative domain-invariant representations. It is important to note that the generator network weights are updated twice during the adversarial training in first and last step respectively.

4. EXPERIMENTAL SETUP

4.1. Dataset

The ADOS-2 dataset is composed of semi-structured activities involving a child and an interlocutor, who is trained to examine behaviours related to ASD. A typical ADOS-2 session lasts between 40-60 minutes and consists of varying subtasks designed to elicit responses from the child under different social and interactive circumstances. In this work, we look at administrations of Module-3 which are intended for verbally-fluent children. Further, we restrict to the *Emotions* and *Social Difficulties & Annoyance* subtasks since they elicit spontaneous speech from the child under significant cognitive load. In the *Emotions* subtask the child is asked to recognize different objects that trigger various emotions within them and share their perceptions on the same. The *Social Difficulties & Annoyance* subtask explores the child’s thoughts regarding various social problems faced at home or school. The dataset consists of recordings from 165 children (86 ASD, 79 Non-ASD) collected from two different clinical centers: University of Michigan Autism and Communication Disorders Center (UMACC) and Cincinnati Children’s Medical Center (CCHMC). Further details are presented in Table 1.

4.2. Features and neural network architecture

In all experiments we used 23-dimensional MFCC features with mean and variance normalized at session level. The features were extracted using the Kaldi¹ toolkit with a frame-

length of 40ms and frame-shift of 20ms. Features were spliced with a context of 15 frames yielding a sample of dimension 31×23 . Consecutive samples were chosen with an interval of 15 frames in order to minimize overlap during DNN training.

The generator $G(\cdot)$ consists of a bidirectional long short term memory (BLSTM) layer followed by four dense layers consisting of 128, 64, 16 and 16 neurons respectively. Certain settings have smaller training data compared to others, hence the number of parameters were reduced to prevent over-fitting. The *speaker classifier* $C(\cdot)$ consists of two dense layers with 16 neurons each, while the *domain discriminator* $D(\cdot)$ consists of one dense layer of 16 neurons. Rectified linear units (ReLU) layers were used as activation functions for all the layers, and both dropout ($p = 0.2$) and batch normalization were applied to every hidden layer for regularization.

4.3. Baselines

We have compared the performance of our systems with two systems. The first system (*Pre-Train*) is composed of only the feature generator and the speaker classifier blocks. This system is trained with source data and directly tested on target data, the goal being to check whether domain adversarial training provides any improvement over pre-training. The second model uses the same architecture, except the training data is augmented with target domain data. Since target labels are not available during domain adversarial training, this system (*Upper-Bound*) serves as an upper bound for the performance.

4.4. Cross-Domain Design

To address the variability resulting from child age and location differences, we designed two sets of experiments: First, we partitioned the data according to age groups and chose the two farthest groups from both locations as the source and the target data. In (*Exp 1*), we selected sessions of kids (≥ 10 yrs) as the source and sessions of kids (≤ 5 yrs) as target data. Later in (*Exp 2*), we reversed the source and target data and repeated the same experiment to address domain shift in the other direction.

Second, we divided the sessions based on their locations. To control for variability sources, we further divided the sessions from each location into 3 age groups of (≤ 5 yrs, 5-10 yrs, ≥ 10 yrs) and conducted separate experiments within each group. In (*Exp 3*), for each age group we considered recordings from Cincinnati as source data and recordings from Michigan as target data. Later, in *Exp 4* we reversed the source and target data and conducted the same experiment).

We check for complementary information in embeddings extracted from GAN and GR using score fusion and embedding fusion. For the score fusion system, we estimate class distribution for a test sample by computing posterior means from GAN and GR models. For the embedding fusion system, we extract embeddings from the output of the generator block for both source and target data for GAN and GR. We then concatenate GAN and GR embeddings and train a separate

¹<https://github.com/kaldi-asr/kaldi>

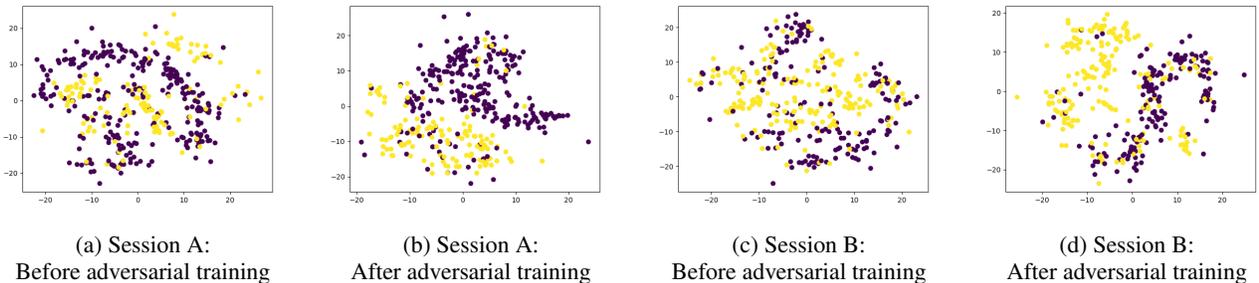


Fig. 3: TSNE plots of the most discriminative 2 components of the generator output corresponding to the classes

Table 2: Mean F1-score (%) treating child age as domain shift

Systems	Exp 1(%)	Exp 2(%)
Pre-Train	73.40	63.69
GAN	78.27	71.21
GR	78.53	72.26
Score Fusion	78.86	71.61
Embed. Fusion	78.38	71.95
Upperbound	85.65	86.29

neural network model with similar architecture to the GAN and GR models, using the source data. Finally, the fused embeddings of the target data are fed to the trained network to check classification performance.

For all experiments, we update model weights using Adam optimizer ($\text{lr} = 0.001$, $\beta_1 = 0.9$, $\beta_2 = 0.999$, $\epsilon = 1e-8$) to minimize categorical cross-entropy loss. Accuracy on a set of held-out sessions from the source corpora is used for early stopping during both pre-training and domain adversarial training. During evaluation, we discard the domain discriminator part. The 23-dimensional features from the audio session are fed to the network consisting of the *generator* $G(\cdot)$ and the *speaker classifier* $C(\cdot)$ to estimate speaker labels at sample-level. Since many sessions in our corpus contain imbalanced class distributions (more samples from adult than child), we estimate classification performance using the mean (unweighted) F1-score.

5. RESULTS AND ANALYSIS

From Tables 2 and 3, we observe that both GAN and GR outperform the baselines in age-based and location-based experiments. In general, GR performs slightly better than GAN in a majority of settings. Among the age-based experiments, we observe that Exp 2 which consists of kids aged ≥ 10 yrs as target data, degrades in accuracy for all models. A possible reason is that older kids with well-developed vocal tract and speaking skills are harder (i.e., more confusable) for the model to discriminate from adult speakers. Interestingly, domain adaptation returns a greater relative improvement over pre-training in Exp 2 (13.45%) than Exp 1 (7.43%).

Among the location-based experiments, the age group ≥ 10 yrs possibly represents the largest domain shift (on the basis of Pre-Train vs Upper-Bound performances). Similar to the age-based experiment, domain adversarial learning

Table 3: Mean F1-score (%) treating collection center as domain shift

Systems	Exp 3(%)			Exp 4(%)		
	≤ 5 yrs	5-10 yrs	≥ 10 yrs	≤ 5 yrs	5-10 yrs	≥ 10 yrs
Pre-Train	79.55	79.23	67.69	82.12	78.16	72.68
GAN	82.14	80.32	73.32	85.03	82.32	76.72
GR	81.74	80.60	73.57	84.53	82.96	76.61
Score Fusion	82.13	80.64	73.46	85.21	83.20	76.85
Embed. Fusion	82.39	80.31	73.19	82.72	82.87	75.33
Upper-bound	87.72	87.56	86.74	90.67	89.47	87.80

returns the largest relative improvement for kids ≥ 10 yrs. Interestingly, improvements in adversarial learning for kids in 5-10 yrs age group are different in Exp 3 and Exp 4. This hints that domain shifts (in this age group) are currently modeled to different extent by GAN and GR, indicating that different modeling techniques should be explored to address this issue. Score fusion performs the best among all the proposed methods, suggesting the presence of complementary information between GAN and GR methods.

As a qualitative analysis, we present TSNE visualizations of the generator outputs for target data from two sessions of Exp 4 in Figure 3. We plot the embeddings before and after GAN training. In both cases, it is evident from the plots that pre-trained embeddings exhibit confusion between child and adult classes, while GAN training increases the discriminative information between them.

6. CONCLUSION

Previous studies have established the potential of adversarial learning for addressing domain mismatch. In this work, we have applied domain adversarial training to enhance the speaker classification performance in autism diagnosis sessions. We have used 2 different methods (*GAN* and *GR*) for learning domain invariant features, and show that domain adversarial training improves the speaker classification performance by a significant margin. Further, we improved the performance further by fusing at the embedding-level and score-level. While our proposed approaches provide improvements over the baseline, the possible upper bound performance implies still significant room for improvement. In the future, we would like to extend adversarial learning to different GAN variants and tasks in the speech pipeline, for example, child ASR.

7. REFERENCES

- [1] Joanne Volden and Catherine Lord, “Neologisms and idiosyncratic language in autistic speakers,” *J. Autism and Developmental Disorders*, vol. 21, no. 2, pp. 109–130.
- [2] So Hyun Kim, Rhea Paul, Helen Tager-Flusberg, and Catherine Lord, *Language and Communication in Autism*, chapter 10, American Cancer Soc., 2014.
- [3] Sabine V Huemer and Virginia Mann, “A comprehensive profile of decoding and comprehension in autism spectrum disorders,” *J. Autism and Developmental Disorders*, vol. 40, no. 4, pp. 485–493, 2010.
- [4] Centers for Disease Control, Prevention (CDC, et al., “Mental health in the united states: parental report of diagnosed autism in children aged 4-17 years—united states, 2003-2004.,” *MMWR. Morbidity and mortality weekly report*, vol. 55, no. 17, pp. 481, 2006.
- [5] Jon Baio et al., “Prevalence of autism spectrum disorder among children aged 8 years—autism and developmental disabilities monitoring network, 11 sites, united states, 2014,” *MMWR Surveillance Summaries*, vol. 67, no. 6, pp. 1, 2018.
- [6] Daniel Bone, Somer L Bishop, Matthew P Black, Matthew S Goodwin, Catherine Lord, and Shrikanth S Narayanan, “Use of machine learning to improve autism screening and diagnostic instruments: effectiveness, efficiency, and multi-instrument fusion,” *J. Child Psychology and Psychiatry*, vol. 57, no. 8, pp. 927–937, 2016.
- [7] Daniel Bone, Somer Bishop, Rahul Gupta, Sungbok Lee, and Shrikanth S Narayanan, “Acoustic-prosodic and turn-taking features in interactions with children with neurodevelopmental disorders.,” in *Interspeech*, 2016, pp. 1185–1189.
- [8] Manoj Kumar, Rahul Gupta, Daniel Bone, Nikolaos Malandrakis, Somer Bishop, and Shrikanth S Narayanan, “Objective language feature analysis in children with neurodevelopmental disorders during autism assessment.,” in *INTERSPEECH*, 2016, pp. 2721–2725.
- [9] Sungbok Lee, Alexandros Potamianos, and Shrikanth Narayanan, “Acoustics of childrens speech: Developmental changes of temporal and spectral parameters,” *The J. Acoustical Soc. of America*, vol. 105, no. 3, pp. 1455–1468, 1999.
- [10] Garrett Wilson and Diane J Cook, “A survey of unsupervised deep domain adaptation,” *arXiv preprint arXiv:1812.02849*, 2019.
- [11] Yaroslav Ganin, Evgeniya Ustinova, Hana Ajakan, Pascal Germain, Hugo Larochelle, François Laviolette, Mario Marchand, and Victor Lempitsky, “Domain-adversarial training of neural networks,” *The J. Machine Learning Research*, vol. 17, no. 1, pp. 2096–2030, 2016.
- [12] Gautam Bhattacharya, Jahangir Alam, and Patrick Kenny, “Adapting end-to-end neural speaker verification to new languages and recording conditions with adversarial training,” in *IEEE Int. Conf. Acoustics, Speech and Signal Process. (ICASSP)*. IEEE, 2019, pp. 6041–6045.
- [13] Maryam Najafian and John HL Hansen, “Speaker independent diarization for child language environment analysis using deep neural networks,” in *2016 IEEE Spoken Lang. Technol. Workshop (SLT)*. IEEE, 2016, pp. 114–120.
- [14] Tianyan Zhou, Weicheng Cai, Xiaoyan Chen, Xiaobing Zou, Shilei Zhang, and Ming Li, “Speaker diarization system for autism children’s real-life audio data,” in *2016 10th Int. Symp. Chinese Spoken Lang. Processing (ISCSLP)*. IEEE, 2016, pp. 1–5.
- [15] L. Sun, J. Du, T. Gao, Y. Lu, Y. Tsao, C. Lee, and N. Ryant, “A novel lstm-based speech preprocessor for speaker diarization in realistic mismatch conditions,” in *2018 IEEE Int. Conf. Acoustics, Speech and Signal Process. (ICASSP)*, April 2018, pp. 5234–5238.
- [16] Alejandrina Cristia, Shobhana Ganesh, Marisa Casillas, and Sriram Ganapathy, “Talker diarization in the wild: The case of child-centered daylong audio-recordings,” in *Interspeech 2018*, 2018, pp. 2583–2587.
- [17] Jiamin Xie, Leibny Paola Garcia-Perera, Daniel Povey, and Sanjeev Khudanpur, “Multi-plda diarization on childrens speech,” in *Interspeech 2019*, 2019, pp. 376–380.
- [18] Mohammed Abdelwahab and Carlos Busso, “Domain adversarial for acoustic emotion recognition,” *IEEE/ACM Trans. Audio, Speech, and Lang. Process.*, vol. 26, no. 12, pp. 2423–2435, 2018.
- [19] Aditay Tripathi, Aanchan Mohan, Saket Anand, and Maneesh Singh, “Adversarial learning of raw speech features for domain invariant speech recognition,” in *2018 IEEE Int. Conf. Acoustics, Speech and Signal Process. (ICASSP)*. IEEE, 2018, pp. 5959–5963.
- [20] Sining Sun, Ching-Feng Yeh, Mei-Yuh Hwang, Mari Ostendorf, and Lei Xie, “Domain adversarial training for accented speech recognition,” in *2018 IEEE Int. Conf. Acoustics, Speech and Signal Process. (ICASSP)*. IEEE, 2018, pp. 4854–4858.
- [21] Pavel Denisov, Ngoc Thang Vu, and Marc Ferras Font, “Unsupervised domain adaptation by adversarial learning for robust speech recognition,” in *Speech Commun.; 13th ITG-Symp.* VDE, 2018, pp. 1–5.
- [22] Gautam Bhattacharya, Joao Monteiro, Jahangir Alam, and Patrick Kenny, “Generative adversarial speaker embedding networks for domain robust end-to-end speaker verification,” in *IEEE Int. Conf. Acoustics, Speech and Signal Process. (ICASSP)*. IEEE, 2019, pp. 6226–6230.