# DEEP IMPLICIT DISTRIBUTION ALIGNMENT NETWORKS FOR CROSS-CORPUS SPEECH EMOTION RECOGNITION

*Yan Zhao, Jincen Wang, Yuan Zong\*, Wenming Zheng\*, Hailun Lian, and Li Zhao*

Key Laboratory of Child Development and Learning Science of Ministry of Education,
Southeast University, Nanjing 210096, China
{zhaoyan, 220222338, xhzongyuan, wenming_zheng, lianhailun, zhaoli}@seu.edu.cn

## ABSTRACT

In this paper, we propose a novel deep transfer learning method called deep implicit distribution alignment networks (DIDAN) to deal with cross-corpus speech emotion recognition (SER) problem, in which the labeled training (source) and unlabeled testing (target) speech signals come from different corpora. Specifically, DIDAN first adopts a simple deep regression network consisting of a set of convolutional and fully connected layers to directly regress the source speech spectrums into the emotional labels such that the proposed DIDAN can own the emotion discriminative ability. Then, such ability is transferred to be also applicable to the target speech samples regardless of corpus variance by resorting to a well-designed regularization term called implicit distribution alignment (IDA). Unlike widely-used maximum mean discrepancy (MMD) and its variants, the proposed IDA absorbs the idea of sample reconstruction to implicitly align the distribution gap, which enables DIDAN to learn both emotion discriminative and corpus invariant features from speech spectrums. To evaluate the proposed DIDAN, extensive cross-corpus SER experiments on widely-used speech emotion corpora are carried out. Experimental results show that the proposed DIDAN can outperform lots of recent state-of-the-art methods in coping with the cross-corpus SER tasks.

*Index Terms*— Cross-corpus speech emotion recognition, speech emotion recognition, deep transfer learning, transfer learning, deep learning.

## 1. INTRODUCTION

As a typical task of affective computing and speech signal processing, research of SER seeks to empower the computers to automatically understand the emotional states, e.g., *Happy*, *Fear*, and *Disgust*, from the speech signals. It has constantly been under the spotlight over past several decades [1, 2] and lots of promising SER methods have been proposed [3, 4, 5]. However, it is noted that numerous interference factors, *e.g.*, language gap, speaker difference, and corpus variance between the training and testing speech signals, still hinder the

---

* indicates the corresponding authors.

possibility of existing well-performing SER methods to move from the laboratory to the practical scenes. This is because that these interference factors would leads to a feature distribution mismatch between the training and testing speech signals and hence remarkably degrade the performance of most well-performing SER methods. To overcome this shortcoming, in recent years some researchers have drawn their attention to a more challenging but fascinating SER issue, *a.k.a.*, cross-corpus SER [6]. Different from the conventional SER, the labeled training and unlabeled testing speech signals in cross-corpus SER belong to different speech emotion corpora. We also refer the training and testing samples/corpora/features/signals as the source and target ones, respectively.

The earliest contribution to cross-corpus SER can be traced to the work of [6], in which Schuller *et al.* proposed a series of feature normalization methods including corpus normalization, speaker normalization, and corpus-speaker normalization to eliminate the corpus difference between the source and target speech signals. Subsequently, several researchers tried to treat cross-corpus SER as a transfer learning task and proposed lots of well-performing transfer subspace learning and deep transfer learning methods. For example, Liu *et al.* [7] proposed a novel transfer subspace learning method called domain-adaptive subspace learning (DoSL) to learn a common subspace to remove the feature distribution mismatch between the source and target speech samples by minimizing their marginal MMD. In the work of [8], Zhao *et al.* proposed a novel deep regression method called deep transductive transfer regression networks (DTTRN), whose major module designed for adapting source and target speech feature distributions are still based on the variant of MMD, *i.e.*, multi-kernel MMD. Besides MMD based methods, adversarial learning is also widely-used in dealing with the cross-corpus SER problem [9, 10]. Unlike MMD and its variants, these methods are not straightforward ones. This means they align the feature distribution gap using an implicit way, *i.e.*, leveraging a domain (corpus) discriminator to disable the networks to be aware of corpus variance.

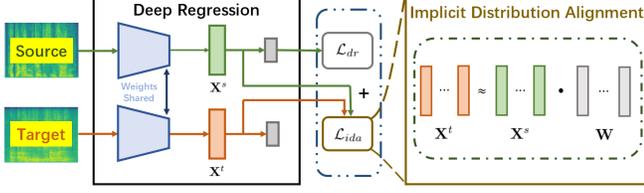Inspired by the success of the adversarial learning based

**Fig. 1.** Overview Structure and Basic Idea of the Proposed DIDAN for Dealing with the Cross-Corpus SER Problem.

methods, in this paper we also focus on the research of cross-corpus SER from the angle of implicit distribution alignment (IDA) and propose a novel deep transfer learning method called deep implicit distribution aligned networks (DIDAN). The major contribution of the proposed DIDAN is designing a novel IDA term to calibrate the feature distribution gap caused by the speech corpus variance. Moreover, different from both MMD and adversarial learning based methods, our IDA performs corpus invariant feature learning by enforcing the learned target speech features to be sparsely reconstructed by the source ones. Hence, all target samples gradually become more involved in the source speech corpus and eventually share the same or similar feature distribution with the source ones.

## 2. PROPOSED METHOD

In this section, we will address the proposed DIDAN for dealing with the cross-corpus SER problem in detail. Suppose we are given a source speech emotion corpus whose samples are denoted by $\mathcal{D}_s = \{(\mathcal{X}_i^s, \mathbf{y}_i^s)\}_{i=1}^{N_s}$, where $\mathcal{X}_i^s$ and $\mathbf{y}_i^s$ are the $i^{th}$ source speech spectrum and its corresponding one-hot emotion class label, and $N_s$ is the source sample number, respectively. Similarly, the unlabeled speech samples from the target corpus can be denoted by $\mathcal{D}_t = \{(\mathcal{X}_i^t, \mathbf{y}_i^t)\}_{i=1}^{N_t}$, where $\mathcal{X}_i^t$ and $N_t$ represent the $i^{th}$ target speech spectrum and the target sample number, respectively. To make the readers better understand the proposed DIDAN, we draw an overall picture shown in Fig. 1 to illustrate the basic idea and network structure. As Fig. 1 shows, our DIDAN has two major parts including **Deep Regression** and **Implicit Distribution Alignment**. In what follows, we will describe them in sequence.

### 2.1. Deep Regression

In our DIDAN, we first build a simple deep regression consisting of a set of convolutional and fully connected layers to bridge the source speech spectrums and their corresponding emotion labels to own the emotion discriminative ability. To achieve this goal, we can optimize the deep regression loss as follows:

$$\mathcal{L}_{dr} = \frac{1}{N_s} \sum_{i=1}^{N_s} J(g(f(\mathcal{X}_i^s)), \mathbf{y}_i^s), \qquad (1)$$

where $J(\cdot)$, $f$ and $g$ denote the cross-entropy loss function, convolution and full connection operations, respectively. It is clear to see that by feeding the source speech samples to the deep regression network and minimizing the above loss function, the DIDAN can gradually be aware of how to distinguish different emotional speech signals.

### 2.2. Implicit Distribution Alignment

Subsequently, we design a novel loss function called IDA to enable DIDAN to be also applicable to recognizing the emotions of target speech signals. Specifically, instead of measuring and narrowing the feature distribution gap like MMD based methods, we would like to make each target speech feature learned by DIDAN possibly look like a source one. To this end, together with deep regression loss function, the following regularization term should be also included in optimization of DIDAN, which can be expressed as:

$$\mathcal{L}_{ida} = \|\mathbf{X}^t - \mathbf{X}^s \mathbf{W}\|_F^2 + \alpha \|\mathbf{W}\|_1, \qquad (2)$$

where $\mathbf{X}^s = [f(\mathcal{X}_1^s), \cdots, f(\mathcal{X}_{N_s}^s)] \in \mathbb{R}^{d \times N_s}$, $\mathbf{X}^t = [f(\mathcal{X}_1^t), \cdots, f(\mathcal{X}_{N_t}^t)] \in \mathbb{R}^{d \times N_t}$, $\mathbf{W} = [\mathbf{w}_1, \cdots, \mathbf{w}_{N_t}] \in \mathbb{R}^{N_s \times N_t}$ is a reconstruction coefficient matrix whose $i^{th}$ column corresponds to the $i^{th}$ target speech sample, and $\alpha$ is a trade-off parameter. It is also noted that $\|\mathbf{W}\|_1 = \sum_{i=1}^{N_t} \|\mathbf{w}_i\|_1$ is a $L_1$ norm with respect to the reconstruction coefficient matrix. By minimizing such norm, DIDAN would produce a sparse $\mathbf{w}_i$, which means only a few source samples are needed to reconstruct $i^{th}$ target one.

### 2.3. Total Loss Function

By combining Eqs.(3) and (2), we will arrive at the final total loss function for learning DIDAN, whose corresponding optimization problem is as follows:

$$\min_{\theta_f, \theta_g, \mathbf{W}} \mathcal{L}_{dr} + \lambda \mathcal{L}_{ida}, \qquad (3)$$

where $\theta_f$ and $\theta_g$ denote the network parameters corresponding to the convolutional operation $f$ and full connection operation $g$, respectively, and $\lambda$ is the trade-off parameter balancing the deep regression and IDA losses.

## 3. EXPERIMENTS

### 3.1. Speech Emotion Corpora and Experimental Setup

To evaluate the proposed DIDAN, three public available speech emotion corpora, *i.e.*, EmoDB (B) [11], eNTERFACE (E) [12], and CASIA (C) [13], are employed to design the cross-corpus SER experiments. EmoDB is a German speech emotion corpus consisting of 535 speech samples from 10 speakers in total. These speakers were requested to perform the seven pre-defined emotional contexts including *Happy*

**Table 1**. The sample statistics of corpora used in the designed six cross-corpus SER tasks.

| Tasks | Speech Corpus (# Samples of Each Emotion) | Total |
|---|---|---|
| B → E | B (AN: 127, SA: 62, FE: 69, HA: 71, DI: 46) | 375 |
| E → B | E (AN: 211, SA: 211, FE: 211, HA: 208, DI: 211) | 1052 |
| B → C | B (AN: 127, SA: 62, FE: 69, HA: 71, NE: 79) | 408 |
| C → B | C (AN: 200, SA: 200, FE: 200, HA: 200, NE: 200) | 1000 |
| E → C | E (AN: 211, SA: 211, FE: 211, HA: 208, SU: 211) | 1052 |
| C → E | C (AN: 200, SA: 200, FE: 200, HA: 200, SU: 200) | 1000 |

(HA), *Sad* (SA), *Disgust* (DI), *Angry* (AN), *Fear* (FE), *Neutral* (NE), and *Boredom*. Different from EmoDB, eNTER-FACE is an English audio-visual bimodal emotion database and hence in the experiments we only adopt its audio data. It is collected from 43 individuals resulting 1582 samples, each of which is labeled as one of six basic emotions including HA, SA, FE, AN, *Surprise* (SU) and *Disgust* (DI). As for CASIA, it is a Chinese speech emotion corpus. It has four different speakers and totally 1200 speech samples. Each speaker is required to perform utterances with six different emotions, *i.e.*, AN, SA, FE, HA, NE, and SU, respectively.

In the task of cross-corpus SER, one speech corpus is served as the source one and the other different corpus as the target one. Therefore, by alternatively choosing either two of the above three speech emotion corpora and meanwhile extracting the speech samples sharing the same emotion labels, we are able to obtain six cross-corpus SER tasks, which can be denoted by B → E, E → B, B → C, C → B, E → C, and C → E, respectively. Note that the right and left corpora of the arrow correspond to the source and target ones. Table 1 summarizes the statistical information of speech samples associated with these six tasks. Moreover, following the pioneer work in cross-corpus SER [6], we adopt unweighted average recall (UAR), which is defined as the average of the prediction accuracy per class, to serve as the performance metric.

### 3.2. Comparison Methods and Implementation Detail

In order to evaluate the effectiveness of the proposed DIDAN, several state-of-the-art transfer subspace learning and deep transfer learning methods are chosen to conduct the comparison experiments. The transfer subspace learning methods include transfer component analysis (TCA) [14], geodesic flow kernel (GFK) [15], subspace alignment (SA) [16], domain-adaptive subspace learning (DoSL) [7], joint distribution adaptive regression (JDAR) [17] and joint distribution implicitly aligned subspace learning (JIASL) [18]. Note that for these subspace learning methods, we use openSMILE toolkit [19] to extract IS09 speech feature sets [20], which consists of 32 low-level acoustic descriptors and 12 statistical functions, to describe the speech signals in all three speech corpora. In the experiments, the elements in IS09 feature set are normalized between 0 and 1. As for the deep learning ones, deep adaptation network (DAN) [21], domain-adversarial neural network (DANN) [22], and conditional domain adversarial network (CDAN) [23], and deep subdomain adaptation network (DSAN) [24] are adopted. Unlike the subspace learning methods, the inputs of the deep neural networks are the speech spectrums converted by applying Fourier Transformation to the original speech signals instead of the hand-crafted IS09 feature set.

In the experiments, we follow existing cross-corpus SER works by searching the hyper-parameters for all the comparison methods from a preset parameter interval to report their best results. Specifically, for TCA, GFK, and SA, we search its reduced feature dimension from $[5 : 5 : d_{max}]$. As for DoSL and JDAR, both of them have two trade-off parameters, *i.e.*, $\lambda$ controlling the distribution alignment term and $\mu$ corresponding to sparsity of projection matrix, whose searching interval is set as $[5 : 5 : 200]$. For all the deep learning methods, VGG-11 is adopted to serve as the backbone and hence the speech spectrums are resized to $224 \times 224$ pixels. We also include the original VGG-11 in the comparison. We search the trade-off parameters of deep learning methods from the hyper-parameter set $\{0.1 : 0.1 : 1, 5, 10, 50, 100\}$. The mini-batch stochastic gradient descent strategy is used for learning the optimal parameters of the deep learning methods. The batch sizes of the source and target speech samples are both fixed at 32.

### 3.3. Results and Discussions

The detailed experimental results are given in Table 2. Several interesting observations and conclusions can be obtained.

(1) It is clear to see that the proposed DIDAN method achieved the best average UAR reaching 39.8% among all the methods. Moreover, we also observed that our DIDAN outperformed all the comparison methods including subspace learning and deep learning ones in three of all six cross-corpus SER tasks, *i.e.*, E → B (46.0%), B → C (39.1%), and C → B (54.5%). Nevertheless, it can be found that the performance of our DIDAN is actually very competitive against the best-performing ones in the rest three tasks, *i.e.*, 36.1% (DIDAN) *v.s.* 36.9% (JIASL) in B → E, 31.9% (DIDAN) *v.s.* 32.8% (GFK) in E → C, and 30.9% (DIDAN) *v.s.* 33.2% (JIASL) in C → E. The above observations demonstrated the effectiveness and superior performance of the proposed DIDAN in dealing with the problem of cross-corpus SER.

(2) Overall speaking, the deep learning methods showed more promising performance in dealing with cross-corpus SER tasks than the subspace learning ones, which can be seen from the comparison between their average UAR. Despite of this, it is noticed that some subspace learning methods can still obtain more satisfactory results compared with the deep learning ones when coping with several cross-corpus SER tasks, *e.g.*, JIASL (36.9%) in B → E and GFK (32.8%) in E → C. This may be because that the hand-crafted speech feature set used in subspace learning methods, *i.e.*, IS09, is more discriminative and corpus invariant than the deep

**Table 2**. The results of all transfer learning methods for cross-corpus SER tasks, where the best results are highlighted in bold.

| Method | | B→E | E→B | B→C | C→B | E→C | C→E | Average |
|---|---|---|---|---|---|---|---|---|
| Subspace Learning | SVM | 28.9 | 23.6 | 29.6 | 35.0 | 26.1 | 25.1 | 28.1 |
| | TCA | 30.5 | 44.0 | 33.4 | 45.1 | 31.1 | 32.3 | 36.1 |
| | GFK | 32.1 | 42.5 | 33.1 | 48.1 | **32.8** | 28.1 | 36.1 |
| | SA | 33.5 | 43.9 | 35.8 | 44.0 | 32.6 | 28.2 | 36.3 |
| | DoSL | 36.1 | 39.0 | 34.4 | 45.8 | 30.4 | 31.6 | 36.2 |
| | JDAR | 36.3 | 40.0 | 31.1 | 46.3 | 32.4 | 31.5 | 36.3 |
| | JIASL | **36.9** | 44.1 | 36.5 | 49.3 | 30.5 | **33.2** | 38.4 |
| Deep Learning | VGG-11 | 32.8 | 38.8 | 36.4 | 50.0 | 27.1 | 30.0 | 35.9 |
| | DAN | 35.2 | 39.2 | 36.7 | 51.6 | 28.5 | 32.5 | 37.3 |
| | JAN | 34.9 | 39.6 | 37.4 | 52.1 | 27.9 | 29.8 | 37.0 |
| | DANN | 35.0 | 43.6 | 37.6 | 52.3 | 28.9 | 30.0 | 37.9 |
| | CDAN | 32.9 | 40.9 | 37.9 | 49.5 | 30.7 | 30.5 | 37.1 |
| | DSAN | 35.6 | 44.0 | 38.5 | 53.4 | 30.3 | 31.7 | 38.9 |
| | DIDAN (Ours) | 36.1 | **46.0** | **39.1** | **54.5** | 31.9 | 30.9 | **39.8** |

**Table 3**. The ablation analysis for the proposed DIDAN.

| Method | B→E | E→B | B→C |
|---|---|---|---|
| VGG-11 (DIDAN *w/o* IDA) | 32.8 | 38.8 | 36.4 |
| DIDAN *w* nonSR-IDA | 34.5 | 42.9 | 37.7 |
| **DIDAN *w* IDA** | **36.1** | **46.0** | **39.1** |
| DAN (MMD) | 35.2 | 39.2 | 36.7 |
| DANN (Adversarial Learning) | 35.0 | 43.6 | 37.7 |

features directly learned from the speech spectrums by the backbone (VGG-11) used in deep learning ones for these tasks. By further comparing the results of all the methods in these tasks, it is clear that the deep learning methods mostly performed poorer compared with the subspace learning ones, which supports our explanations and analysis.

(3) It can be observed that nearly all the methods cannot well cope with the cross-corpus SER tasks between eNTER-FACE and CASIA. We believe that this may attribute to the large differences between these two speech corpora. According to the works of [12, 13], it is well known that eNTER-FACE is an English speech corpus whose samples are elicited by well-designed induced paradigm, while CASIA is a Chinese one, in which all speakers are requested to simulate different emotions.

### 3.4. Going Deeper into IDA of DIDAN

As described previously, one of the major contributions in our DIDAN is the IDA loss, which aims to improve the robustness of DIDAN to the corpus variance by enforcing the target speech features to be sparsely reconstructed by the learned source ones. To see whether IDA indeed works, we select three cross-corpus SER tasks, *i.e.*, B → E, E → B, and B → C, as representatives to conduct additional experiments. Besides the original DIDAN (denoted by DIDAN *w* IDA), another four methods are chosen including VGG-11 (DIDAN *w/o* IDA), DIDAN without sparsity regularization term (denoted by DIDAN *w* nonSR-IDA), DAN (MMD), and DANN

(Adversarial Learning). Experimental results are depicted in Table 3. From the results in first three rows of Table 3, it can be concluded that the proposed implicit distribution alignment method can remarkably improve the corpus invariant ability of the deep regression model. In addition, we can also observe that our DIDAN outperformed DAN and DANN, which adopt MMD and GAN to align the feature distribution gap, respectively. This verifies the superiority of the proposed IDA in DIDAN over these two widely-used strategies for distribution alignment.

## 4. CONCLUSION

In this paper, we have presented a novel deep transfer learning method called DIDAN for dealing with the problem of cross-corpus SER. Unlike most of existing deep learning methods, our DIDAN removes the feature distribution mismatch between the source and target speech signals with an implicit manner, *i.e.*, enforcing the target deep features to be sparsely reconstructed by the source ones. Hence, the deep regression model only supervised by the source label information would be able to effectively recognize the emotions of unlabeled target speech signals. Extensive experiments were carried out on three widely-used speech emotion corpora to evaluate the performance of the proposed DIDAN. The results showed that compared with recent state-of-the-art transfer subspace learning and deep transfer learning methods, our DIDAN has more promising performance in dealing with cross-corpus SER tasks.

## 5. ACKNOWLEDGEMENTS

## 6. REFERENCES

[1] Moataz El Ayadi, Mohamed S. Kamel, and Fakhri Karray, "Survey on speech emotion recognition: Features, classification schemes, and databases," *Pattern Recognition*, vol. 44, no. 3, pp. 572–587, 2011.

[2] Björn W Schuller, "Speech emotion recognition: Two decades in a nutshell, benchmarks, and ongoing trends," *Communications of the ACM*, vol. 61, no. 5, pp. 90–99, 2018.

[3] Yuan Zong, Wenming Zheng, Zhen Cui, and Qiang Li, "Double sparse learning model for speech emotion recognition," *Electronics Letters*, vol. 52, no. 16, pp. 1410–1412, 2016.

[4] Seyedmahdad Mirsamadi, Emad Barsoum, and Cha Zhang, "Automatic speech emotion recognition using recurrent neural networks with local attention," in *ICASSP*. IEEE, 2017, pp. 2227–2231.

[5] Cheng Lu, Yuan Zong, Wenming Zheng, Yang Li, Chuangao Tang, and Björn W Schuller, "Domain invariant feature learning for speaker-independent speech emotion recognition," *IEEE/ACM Transactions on Audio, Speech, and Language Processing*, vol. 30, pp. 2217–2230, 2022.

[6] Bjorn Schuller, Bogdan Vlasenko, Florian Eyben, Martin Wöllmer, Andre Stuhlsatz, Andreas Wendemuth, and Gerhard Rigoll, "Cross-corpus acoustic emotion recognition: Variances and strategies," *IEEE Transactions on Affective Computing*, vol. 1, no. 2, pp. 119–131, 2010.

[7] Na Liu, Yuan Zong, Baofeng Zhang, Li Liu, Jie Chen, Guoying Zhao, and Junchao Zhu, "Unsupervised cross-corpus speech emotion recognition using domain-adaptive subspace learning," in *ICASSP*. IEEE, 2018, pp. 5144–5148.

[8] Yan Zhao, Jincen Wang, Ru Ye, Yuan Zong, Wenming Zheng, and Li Zhao, "Deep transductive transfer regression network for cross-corpus speech emotion recognition," *INTERSPEECH*, pp. 18–22, 2022.

[9] Mohammed Abdelwahab and Carlos Busso, "Domain adversarial for acoustic emotion recognition," *IEEE/ACM Transactions on Audio, Speech, and Language Processing*, vol. 26, no. 12, pp. 2423–2435, 2018.

[10] Bo-Hao Su and Chi-Chun Lee, "Unsupervised cross-corpus speech emotion recognition using a multi-source cycle-gan," *IEEE Transactions on Affective Computing*, 2022.

[11] Felix Burkhardt, Astrid Paeschke, M. Rolfes, Walter F. Sendlmeier, and Benjamin Weiss, "A database of german emotional speech," in *INTERSPEECH*, 2005, pp. 1517–1520.

[12] Olivier Martin, Irene Kotsia, Benoît Macq, and Ioannis Pitas, "The enterface'05 audio-visual emotion database," in *ICDE Workshops*, 2006, p. 8.

[13] Jianhua Tao, Fangzhou Liu, Meng Zhang, and Huibin Jia, "Design of speech corpus for mandarin text to speech," in *The Blizzard Challenge 2008 Workshop*, 2008, pp. 1–4.

[14] Sinno Jialin Pan, Ivor W Tsang, James T Kwok, and Qiang Yang, "Domain adaptation via transfer component analysis," *IEEE Transactions on Neural Networks*, vol. 22, no. 2, pp. 199–210, 2010.

[15] Boqing Gong, Yuan Shi, Fei Sha, and Kristen Grauman, "Geodesic flow kernel for unsupervised domain adaptation," in *CVPR*. IEEE, 2012, pp. 2066–2073.

[16] Basura Fernando, Amaury Habrard, Marc Sebban, and Tinne Tuytelaars, "Unsupervised visual domain adaptation using subspace alignment," in *ICCV*, 2013, pp. 2960–2967.

[17] Jiacheng Zhang, Lin Jiang, Yuan Zong, Wenming Zheng, and Li Zhao, "Cross-corpus speech emotion recognition using joint distribution adaptive regression," in *ICASSP*. IEEE, 2021, pp. 3790–3794.

[18] Cheng Lu, Yuan Zong, Chuangao Tang, Hailun Lian, Hongli Chang, Jie Zhu, Sunan Li, and Yan Zhao, "Implicitly aligning joint distributions for cross-corpus speech emotion recognition," *Electronics*, vol. 11, no. 17, pp. 2745, 2022.

[19] Florian Eyben, Martin Wöllmer, and Björn Schuller, "Opensmile: the munich versatile and fast open-source audio feature extractor," in *ACM Multimedia*, 2010, pp. 1459–1462.

[20] Björn Schuller, Stefan Steidl, and Anton Batliner, "The interspeech 2009 emotion challenge," 2009.

[21] Mingsheng Long, Yue Cao, Jianmin Wang, and Michael Jordan, "Learning transferable features with deep adaptation networks," in *ICML*. PMLR, 2015, pp. 97–105.

[22] Hana Ajakan, Pascal Germain, Hugo Larochelle, François Laviolette, and Mario Marchand, "Domain-adversarial neural networks," *arXiv preprint arXiv:1412.4446*, 2014.

[23] Mingsheng Long, Zhangjie Cao, Jianmin Wang, and Michael I Jordan, "Conditional adversarial domain adaptation," *NIPS*, vol. 31, 2018.

[24] Yongchun Zhu, Fuzhen Zhuang, Jindong Wang, Guolin Ke, Jingwu Chen, Jiang Bian, Hui Xiong, and Qing He, "Deep subdomain adaptation network for image classification," *IEEE Transactions on Neural Networks and Learning Systems*, vol. 32, no. 4, pp. 1713–1722, 2020.