

Unsupervised Cross-Lingual Speech Emotion Recognition Using Domain Adversarial Neural Network

Xiong Cai¹, Zhiyong Wu^{1,2}, Kuo Zhong¹, Bin Su¹, Dongyang Dai¹, Helen Meng^{1,2}

¹Tsinghua-CUHK Joint Research Center for Media Sciences, Technologies and Systems, Shenzhen International Graduate School, Tsinghua University, Shenzhen, China

²Department of Systems Engineering and Engineering Management,

The Chinese University of Hong Kong, Shatin, N.T., Hong Kong SAR, China

{cai-x18, zhongk17, sub18, ddy17}@mails.tsinghua.edu.cn, {zywu, hmmeng}@se.cuhk.edu.hk

Abstract

By using deep learning approaches, Speech Emotion Recognition (SER) on a single domain has achieved many excellent results. However, cross-domain SER is still a challenging task due to the distribution shift between source and target domains. In this work, we propose a Domain Adversarial Neural Network (DANN) based approach to mitigate this distribution shift problem for cross-lingual SER. Specifically, we add a language classifier and gradient reversal layer after the feature extractor to force the learned representation both language-independent and emotion-meaningful. Our method is unsupervised, i. e., labels on target language are not required, which makes it easier to apply our method to other languages. Experimental results show the proposed method provides an average absolute improvement of 3.91% over the baseline system for arousal and valence classification task. Furthermore, we find that batch normalization is beneficial to the performance gain of DANN. Therefore we also explore the effect of different ways of data combination for batch normalization.

Index Terms: speech emotion recognition, domain adversarial learning, cross-lingual, affective representation learning

1. Introduction

With the extensive application of Artificial Intelligence (AI) products in our daily lives, it has become increasingly imperative to design a smarter Human-Computer Speech Interaction (HCSI) system. Speech Emotion Recognition (SER), which aims to infer the emotional state of a speaker from his or her speech [1], has been regarded as a crucial component for a more intelligent HCSI system. Existing SER models [2–4] have achieved satisfactory level results when the training and test data are from the same corpus. However, it is still intractable to build a more robust cross-lingual SER system because of the domain shift between corpora of different languages [5].

Numerous approaches have been proposed to reduce the domain shift problem for cross-corpus or cross-lingual SER. [6] proposes a fine-grained adversarial domain adaptation scheme, which reduces the distribution shift of the same emotion class in different corpora. [7] shows that fine-tuning can effectively improve the recognition results. These methods are promising, but additional labeled data are required, which might not be available since their collection is expensive.

A more practical solution is unsupervised domain adaptation which only demands unlabeled data from related domains. A number of previous studies have explored statistical-based methods to reduce mismatch between domains [8–12]. Specifically, [8, 9] deploy different level of feature normaliza-

tion strategies to minimize the speaker-and-corpus-related effects; [10–12] apply the Maximum Mean Discrepancy (MMD) or Kernel Canonical Correlation Analysis (KCCA) approaches to increase the similarity or correlation of different domains. All these methods reduce the domain shift directly on the original input feature space or its linear transformation space, so the capacity of shift-reduction might be limited. Some other studies [13–15] use variants of autoencoder to learn a concise and common feature representation by incorporating the prior knowledge from unlabeled data into learning. Since the optimization of the autoencoder and emotion classifier is not performed simultaneously, it is not clear whether compressed representations preserve all the emotion information of speech.

Recently, Adversarial Learning (AL), such as Generative Adversarial Network (GAN) [16] and Domain Adversarial Neural Network (DANN) [17], has become an increasingly popular approach for domain adaptation. [18] proposes a GAN-based model for cross-lingual SER and demonstrates significant improvements even for the non-mainstream Urdu language. [19, 20] use a DANN-based framework to learn a speaker-independent representation and greatly improve the single-corpus results. [21] explores the advantage of DANN for cross-corpus SER on three English corpora. Besides, the DANN techniques have also been widely applied in other speech applications such as automatic speech recognition [22] and speaker recognition [23] to deal with the domain mismatch problem.

Inspired by the success of DANN in domain adaptation tasks, this paper proposes a DANN-based approach to reduce the distribution shift for cross-lingual SER. Specifically, based on the primary emotion classification task, a language classifier with Gradient Reversal Layer (GRL) is added to the model as an auxiliary task to help learn language-independent representations. Unlike the studies mentioned above, our approach reduces the distribution shift in a compressed feature space instead of the original input space, and all the modules are trained jointly rather than separately, which makes the model learn emotion-discriminative and language-independent representations more efficiently. Our contribution is two-fold: First, we introduce the DANN framework for cross-lingual SER and achieve significant performance improvements. Second, our study presents that batch normalization (BN) [24] can contribute to improve DANN and explores four different ways of combining data for BN.

The rest of this paper is organized as follows. The proposed method is described in Section 2. The databases and classification scheme are detailed in Section 3. Experimental setup and results analysis are presented in Section 4. Section 5 finalizes the study with conclusions and future directions.

2. Methodology

Firstly, we formulate our cross-lingual SER as the following domain adaptation task. We have a source language corpus with emotion labels as source domain, $D_s = \{(\mathbf{x}_i^s, \mathbf{y}_i^s)\}_{i=1}^n$, and a target language corpus without emotion labels as target domain, $D_t = \{\mathbf{x}_i^t\}_{i=1}^m$, where \mathbf{x}_i^s and $\mathbf{x}_i^t \in \mathbb{R}^{k \times d}$, $\mathbf{y}_i^s \in \{0, 1\}^c$; k , d and c are the number of frames in an utterance, the dimension of each feature frame and the number of emotion categories; n and m are the number of samples in D_s and D_t . Our goal is to learn a reliable emotion classifier from the labeled D_s and the unlabeled D_t , which can be generalized well in D_t .

2.1. Model structure

As shown in Figure 1, the proposed model consists of three modules: encoder (G_f), emotion classifier (G_e) and language classifier (G_l). The encoder structure is mainly adopted from [7], except that we add a batch normalization (BN) layer for the stability of training a DANN model. A 1D convolution layer with ReLU activation takes Mel spectral features as input to capture emotion-related patterns. Then, a max pooling layer with a large stride follows to select the most salient features. Next, an attentive vector $\mathbf{a}_{\hat{f}}$ is extracted from the outputs of max pooling by the following attention formulas:

$$s_i = \frac{\exp(\mathbf{v}^T \hat{\mathbf{f}}_i)}{\sum_j \exp(\mathbf{v}^T \hat{\mathbf{f}}_j)} \quad (1)$$

$$\mathbf{a}_{\hat{f}} = \sum_i s_i \hat{\mathbf{f}}_i \quad (2)$$

where $\hat{\mathbf{f}}_i$ is the i -th feature vector of the output $\hat{\mathbf{f}}$ of max pooling layer and \mathbf{v} is a trainable vector as a global attention query. The motivation behind using this attention mechanism is that emotion-related information is distributed differently over the utterance. This global attention query \mathbf{v} can be used to learn to capture these important emotion patterns. Finally, the attention vector $\mathbf{a}_{\hat{f}}$ is appended to the end of the output $\hat{\mathbf{f}}$ of max pooling along the time dimension, and then all these feature vectors are flattened into a fixed-length vector as the input of the following BN layer. As the final representation, the output of the BN layer \mathbf{f} is fed into emotion classifier (only source domain data) and language classifier (both source and target domain data). As for classifiers, a single dense layer with softmax activation and two output units are used for both classifiers. Besides, a Gradient Reversal Layer (GRL) is inserted between BN layer and language classifier to achieve the goal of adversarial training.

2.2. Adversarial training

DANN [17] is a fairly elegant neural network framework for unsupervised domain adaptation, where unlabeled target domain data can be efficiently utilized to reduce the variations between the source and target domains. Specifically, there are two tasks: a primary target task (e. g., emotion classification) and an auxiliary domain classification task (e. g., language classification). Both tasks share the feature extractor and a GRL is introduced between feature extractor and domain classifier. GRL is a layer without trainable parameters and works as a ‘‘pseudo-function’’ $R(\mathbf{x})$ defined as the following formulas:

$$R(\mathbf{x}) = \mathbf{x} \quad (3)$$

$$\frac{dR}{d\mathbf{x}} = -\beta \mathbf{I} \quad (4)$$

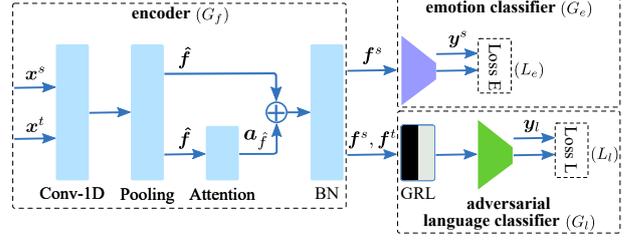


Figure 1: The proposed model structure. BN and GRL are Batch Normalization, and Gradient Reversal Layer respectively.

where \mathbf{I} is an identity matrix and β is a hyper parameter controlling the scale of reversal gradient signal. By using GRL, the trainable parameters before and after GRL are updated in the opposite direction, namely, *adversarial training*. As for the classifiers of the two tasks, parameters are updated to minimize their respective errors. As for the feature extractor, parameters are updated to minimize the error of primary task while maximizing the error of domain classification task, where the latter is implemented by GRL. Therefore, the learned feature could be meaningful for the primary task and indistinguishable for the domain classifier. In terms of our cross-lingual SER, the primary task is emotion classification and the auxiliary task is language classification. Our goal is to learn a representation that retains discriminative information for emotion and reduces variations for languages. Therefore, the learned feature extractor and emotion classifier can be directly applied to target language data.

We use the cross-entropy loss as the training objective for both emotion and language classifiers:

$$L_e(\theta_f, \theta_e) = -\frac{1}{n} \sum_{(\mathbf{x}, \mathbf{y}) \in D_s} \mathbf{y}^T \log G_e(G_f(\mathbf{x}; \theta_f); \theta_e) \quad (5)$$

$$L_l(\theta_f, \theta_l) = -\frac{1}{n+m} \sum_{\mathbf{x} \in D} \mathbf{y}_l^T \log G_l(G_f(\mathbf{x}; \theta_f); \theta_l) \quad (6)$$

where $D = D_s \cup D_t$; L_e and L_l are the losses for emotion and language classifiers respectively; θ_f , θ_e and θ_l represent the trainable parameters for G_f , G_e and G_l respectively; \mathbf{y} is the one-hot encoding for emotion labels from source domain and \mathbf{y}_l is the one-hot encoding for language labels distinguishing source and target languages.

Then, we define the total loss as the weighted sum of the above two losses and directly minimize it for training:

$$L(\theta_f, \theta_e, \theta_l) = \alpha \cdot L_e(\theta_f, \theta_e) + (1 - \alpha) \cdot L_l(\theta_f, \theta_l) \quad (7)$$

where α plays a trade-off for the two losses. Due to the existence of GRL, minimizing the total loss L will actually lead to the following way of parameter update:

$$\theta_e \leftarrow \theta_e - \lambda \cdot \alpha \frac{\partial L_e}{\partial \theta_e} \quad (8)$$

$$\theta_l \leftarrow \theta_l - \lambda \cdot (1 - \alpha) \frac{\partial L_l}{\partial \theta_l} \quad (9)$$

$$\theta_f \leftarrow \theta_f - \lambda \cdot \left(\alpha \frac{\partial L_e}{\partial \theta_f} - (1 - \alpha) \cdot \beta \frac{\partial L_l}{\partial \theta_f} \right) \quad (10)$$

Concretely, θ_e and θ_l are updated for minimizing L_e and L_l respectively, and θ_f is updated for minimizing L_e while maximizing L_l simultaneously. λ and β are the learning rate and the gradient reversal scale of GRL. After training, a feature representation rich in emotional information and indistinguishable from languages will be obtained from the encoder output.

3. Databases

3.1. IEMOCAP

IEMOCAP [25] is an audiovisual database of English dyadic conversations performed by ten professional actors. There are two types of conversations: the scripted ones and the improvised ones (given a certain scenario and topic). This corpus contains a total of 10,039 utterances, where audio, video, text and motion-capture recordings are available. The categorical emotion label and 5-point scales on the dimensions valence, arousal, and dominance (1 - low/negative, 5 - high/positive) are annotated by at least 2 raters. In our study, only audio modal data and dimension label of valence and arousal are used.

3.2. RECOLA

RECOLA [26] is a multimodal database of French dyadic conversations. Participants express emotions spontaneously during a collaborative video conference. Four different modal data of audio, video, electrocardiogram (ECG) and electrodermal activity (EDA) are recorded continuously and synchronously. Continuous valence and arousal labels in the range[-1, 1] are measured by 6 annotators at frame level. Since our goal is to predict emotion on utterance level, the mean value across all frames of an utterance and all annotators are calculated as the final label. Freely available 1,308 audio utterances from 23 speakers are used in our study.

3.3. Classification scheme and input features

In this work, we focus on a binary classification task of valence (negative/positive) and arousal (low/high). In order to obtain binary training labels, we use the same annotation mapping scheme as in [7]. For IEMOCAP, the two ranges [1, 2.5] and (2.5, 5] are categorized as low/negative and high/positive respectively. Similarly, the corresponding two ranges are [-1, 0] and (0, 1] for RECOLA. In terms of input features, 26 log-Mel filter-banks are extracted frame-wise from a single utterance with frame size of 25ms and frame shift of 10ms. The log-Mel feature has a fixed length of 750 frames. The shorter one is padded with the minimum for each dimension in an utterance and the longer one is truncated to 750 frames in the middle.

4. Experiments

4.1. Experimental setup

We use the following configurations for model training. 200 filters with kernel size 10 and stride 3 are used for the 1D convolution layer. The size and stride of max pooling are both set to 30. Adam [27] optimizer and exponential decay learning rate with initial rate 1e-3, decay rate 0.93 for every epoch, and final rate 5e-5 are used to optimize parameters. For the regularization, dropout with rate 0.7 as suggested in [28] is used for the output of encoder; l_1 and l_2 regularization with the weight 5e-3 are used for training RECOLA and IEMOCAP respectively. We train the models for 50 epochs with a batch size of 32, and 30% of data from test set is used as the development set for early stopping. The logMel features are normalized with zero mean and unit variance for each database. All experiments are run five times with different random seeds, and the unweighted average recall (UAR) is chosen as our evaluation criterion.

4.2. Experimental results

4.2.1. Performance of the proposed model

In this section, we compare three trained models: our proposed model (*our*), the baseline model (*base*), and the mono-lingual model (*mono*). The *base* and *mono* model use the structure which consists of the same encoder and emotion classifier only as in Figure 1. The *mono* model is trained and tested on the same database, where 70% samples are used for training, 25% for testing and 5% for early stopping. It provides us with an idea about the best achievable results within each database. We use *Rec* and *Iem* to represent the RECOLA and IEMOCAP database, and *Rec2Iem* means training on *Rec* and testing on *Iem* and vice versa.

Table 1 reports UAR (%) results with standard deviations in parentheses for the three models. Comparing the results of *base* and *our*, our proposed model outperforms the baseline model in all experiments and achieves an average improvement of 3.91%. This result presents that the proposed approach can effectively reduce variations between different languages while retaining the information related to emotions. Therefore, the emotion classifier can benefit from the learned language-independent representation to improve results in target language. To illustrate this, we use Principal Component Analysis (PCA) to project the learned feature representation, i. e., the output of encoder, into 2D space.

Table 1: UAR (%) for baseline and proposed method.

model	Rec2Iem		Iem2Rec		average
	arousal	valence	arousal	valence	
base	62.49(2.96)	54.15(0.51)	60.73(0.45)	58.11(0.51)	58.87
our	71.99(0.33)	54.54(0.77)	63.18(0.32)	61.43(1.38)	62.78
mono ¹	75.55(0.78)	63.20(2.23)	66.28(1.71)	62.90(1.12)	66.98

As shown in Figure 2, regarding language labels, feature representations learned by our adversarial training model (Figure 2(b) left) are more evenly mixed and therefore more indistinguishable than the ones learned by baseline model (Figure 2(a) left); while, as for emotion labels, feature representations learned by our proposed model (Figure 2(b) right) are more separable than the ones learned by baseline model (Figure 2(a) right). Besides, in Figure 3, we plot the training curves of the domain classification loss (blue line) and UAR (green line), and the emotion classification UAR (red line) in development set. It can be seen that, at the early stage of training, the domain loss increases and decreases alternately, and the domain UAR changes oppositely than it, which suggests that the adversarial training itself works well. Moreover, the emotion development set UAR has a similar trend with domain loss, which means that the emotion classifier gets better results when the domain classifier has a higher loss, i. e., the more language-indiscriminative the features are, the better performance of emotion classification will be. These visual results further indicate the effectiveness of the proposed method for cross-lingual SER.

Comparing the results of *base* and *mono* in Table 1, the performance of naive cross-lingual SER (baseline) is 8.11% lower on average than the mono-lingual SER. This result consistent with [5,29] indicates that the distribution shift between different languages will seriously harness the predictive ability of SER.

¹*mono* is trained and tested on *Iem* for *Rec2Iem* setup, and *Rec* for *Iem2Rec* setup

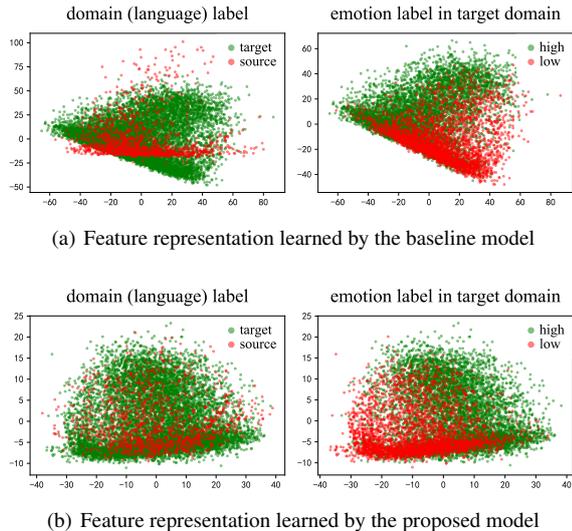


Figure 2: PCA plot of the learned feature representation with language labels (left) and emotion labels (right) for baseline and our proposed model from the “Rec2Iem arousal” training.

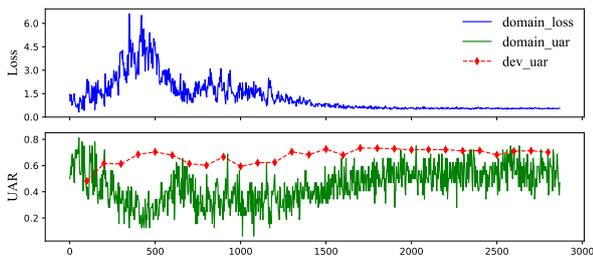


Figure 3: Domain(language) loss, domain UAR, and emotion development set UAR from the “Rec2Iem arousal” training

In addition, another clear conclusion can be obtained that the prediction of arousal is easier than valence regardless of cross-lingual or mono-lingual tasks. Similar results can be found in [5, 30–32]. This is mainly because acoustic features such as energy, pitch and speed are related to arousal [1], but for valence, there is no consensus on how acoustic features correlate with it, and it is more speaker-dependent [28].

4.2.2. Impact of batch normalization on performance

In this section, we first study the effect of BN on widening the performance gap between our DANN-based model and baseline model. As shown in Table 2, our proposed model is only 2.61% higher than baseline on average when BN is not used, which is much lower than the gap of 3.91% as in Table 1 with BN used. A possible reason for this result is that the adversarial training of our model is not as stable as baseline, while BN can help to make the distribution of the representation layer more stable [24]. Therefore, better results can be achieved after using the BN layer in our adversarial training model.

Based on the conclusion reached above, we further explore four different ways of combining data for BN. For the training of DANN model, both source and target data need to be fed into the model. They can be first combined into one mini-batch and then fed into the model, or each of them occupy a min-batch and fed into the model alternatively. For the first data feeding

Table 2: UAR (%) for no batch normalization.

model	Rec2Iem		Iem2Rec		average
	arousal	valence	arousal	valence	
base	58.04(1.90)	52.34(0.79)	58.69(0.95)	53.79(0.87)	55.71
our	63.62(1.74)	52.94(0.36)	59.76(0.77)	57.81(2.07)	58.32

method, three ways of combining data for BN are performed as follows: perform BN on the whole batch, namely **BN1**, which is used for above experiments, where the first half batch (source half) is fed to the emotion classifier (G_e) and the whole batch is fed to the language classifier (G_l); perform BN on the source half batch and whole batch respectively, namely **BN2**; perform BN on the source half batch and target half batch respectively, namely **BN3**. For the second data feeding method, BN is performed on the whole batch from each domain, namely **BN4**.

The evaluation results of the above four types of BN are shown in Table 3. On the one hand, the average results of all **BN1-3** are higher than **BN4**. This proves that it is better to combine data from both source and target domains in one batch than batch them separately. Therefore, it is important for the training of DANN to ensure the guiding gradient signal comes from both source and target domains at each training step. On the other hand, we can also find the average results of both **BN1-2** are better than **BN3**. The main difference between **BN1-2** and **BN3** is whether the input features for G_l is performed BN on the whole batch (**BN1-2**) or on the source and target half separately (**BN3**). This result presents that it is more suitable to feed the language classifier with features performed BN on the entire batch. Besides, it is also worth noting that when training on the smaller database of RECOLA (1,308 utterances), the results of all four settings don’t show significant difference. Therefore, this study empirically suggests that **BN1** or **BN2** is a more recommended way for BN of features, when training the DANN model on a larger corpus.

Table 3: UAR (%) for four different ways of data combination for batch normalization.

model	Rec2Iem		Iem2Rec		average
	arousal	valence	arousal	valence	
BN1	71.99(0.33)	54.54(0.77)	63.18(0.32)	61.43(1.38)	62.78
BN2	72.22(0.18)	54.99(0.84)	62.34(0.92)	61.37(1.15)	62.73
BN3	72.27(0.39)	54.07(0.42)	61.83(1.03)	58.48(0.78)	61.66
BN4	72.01(0.33)	53.68(1.26)	60.95(0.90)	56.77(2.84)	60.85

5. Conclusions

In this paper, we propose a DANN-based approach for cross-lingual SER. Our method works in a completely unsupervised way, where unlabeled target language data is required only. Experimental results show that our method enables the model to focus on the emotion related information, while ignoring the variations between different languages. Moreover, we explore the impact of batch normalization on training DANN models and suggest two practically optimal ways of data combination for batch normalization. For further work, we plan to add more corpora from other languages for the cross-lingual SER task.

6. Acknowledgements

This work is supported by joint research fund of National Natural Science Foundation of China - Research Grant Council of Hong Kong (NSFC-RGC) (61531166002, N_CUHK404/15).

7. References

- [1] M. El Ayadi, M. S. Kamel, and F. Karray, "Survey on speech emotion recognition: Features, classification schemes, and databases," *Pattern Recognition*, vol. 44, no. 3, pp. 572–587, 2011.
- [2] R. Li, Z. Wu, J. Jia, Y. Bu, S. Zhao, and H. Meng, "Towards discriminative representation learning for speech emotion recognition," in *International Joint Conference on Artificial Intelligence (IJCAI)*, 2019, pp. 5060–5066.
- [3] M. A. Jalal, R. K. Moore, and T. Hain, "Spatio-temporal context modelling for speech emotion classification," in *IEEE Automatic Speech Recognition and Understanding Workshop (ASRU)*, 2019, pp. 853–859.
- [4] P. Yenigalla, A. Kumar, S. Tripathi, C. Singh, S. Kar, and J. Vepa, "Speech emotion recognition using spectrogram & phoneme embedding," in *Conference of the International Speech Communication Association (INTERSPEECH)*, 2018, pp. 3688–3692.
- [5] S. M. Feraru, D. Schuller *et al.*, "Cross-language acoustic emotion recognition: An overview and some tendencies," in *International Conference on Affective Computing and Intelligent Interaction (ACII)*, 2015, pp. 125–131.
- [6] H. Zhou and K. Chen, "Transferable positive/negative speech emotion recognition via class-wise adversarial domain adaptation," in *IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, 2019, pp. 3732–3736.
- [7] M. Neumann and N. g. Thang Vu, "Cross-lingual and multilingual speech emotion recognition on English and French," in *IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, 2018, pp. 5769–5773.
- [8] B. Schuller, B. Vlasenko, F. Eyben, M. Wollmer, A. Stuhlsatz, A. Wendemuth, and G. Rigoll, "Cross-corpus acoustic emotion recognition: Variances and strategies," *IEEE Transactions on Affective Computing*, vol. 1, no. 2, pp. 119–131, 2010.
- [9] H. Kaya and A. A. Karpov, "Efficient and effective strategies for cross-corpus acoustic emotion recognition," *Neurocomputing*, vol. 275, pp. 1028–1034, 2018.
- [10] Y. Zong, W. Zheng, T. Zhang, and X. Huang, "Cross-corpus speech emotion recognition based on domain-adaptive least-squares regression," *IEEE Signal Processing Letters*, vol. 23, no. 5, pp. 585–589, 2016.
- [11] P. Song, W. Zheng, S. Ou, X. Zhang, Y. Jin, J. Liu, and Y. Yu, "Cross-corpus speech emotion recognition based on transfer non-negative matrix factorization," *Speech Communication*, vol. 83, pp. 34–41, 2016.
- [12] H. Sagha, J. Deng, M. Gavryukova, J. Han, and B. Schuller, "Cross lingual speech emotion recognition using canonical correlation analysis on principal component subspace," in *IEEE international conference on acoustics, speech and signal processing (ICASSP)*, 2016, pp. 5800–5804.
- [13] J. Deng, Z. Zhang, F. Eyben, and B. Schuller, "Autoencoder-based unsupervised domain adaptation for speech emotion recognition," *IEEE Signal Processing Letters*, vol. 21, no. 9, pp. 1068–1072, 2014.
- [14] J. Deng, R. Xia, Z. Zhang, Y. Liu, and B. Schuller, "Introducing shared-hidden-layer autoencoders for transfer learning and their application in acoustic emotion recognition," in *IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, 2014, pp. 4818–4822.
- [15] J. Deng, X. Xu, Z. Zhang, S. Frühholz, and B. Schuller, "Universum autoencoder-based domain adaptation for speech emotion recognition," *IEEE Signal Processing Letters*, vol. 24, no. 4, pp. 500–504, 2017.
- [16] I. Goodfellow, J. Pouget-Abadie, M. Mirza, B. Xu, D. Warde-Farley, S. Ozair, A. Courville, and Y. Bengio, "Generative adversarial nets," in *Advances in Neural Information Processing Systems*, 2014, pp. 2672–2680.
- [17] Y. Ganin, E. Ustinova, H. Ajakan, P. Germain, H. Larochelle, F. Laviolette, M. Marchand, and V. Lempitsky, "Domain-adversarial training of neural networks," *The Journal of Machine Learning Research*, vol. 17, no. 1, pp. 2096–2030, 2016.
- [18] S. Latif, J. Qadir, and M. Bilal, "Unsupervised adversarial domain adaptation for cross-lingual speech emotion recognition," in *International Conference on Affective Computing and Intelligent Interaction (ACII)*, 2019, pp. 732–737.
- [19] M. Tu, Y. Tang, J. Huang, X. He, and B. Zhou, "Towards adversarial learning of speaker-invariant representation for speech emotion recognition," *arXiv preprint arXiv:1903.09606*, 2019.
- [20] Z. Lian, J. Tao, B. Liu, and J. Huang, "Domain adversarial learning for emotion recognition," *arXiv preprint arXiv:1910.13807*, 2019.
- [21] M. Abdelwahab and C. Busso, "Domain adversarial for acoustic emotion recognition," *IEEE/ACM Transactions on Audio, Speech, and Language Processing*, vol. 26, no. 12, pp. 2423–2435, 2018.
- [22] G. Saon, G. Kurata, T. Sercu, K. Audhkhasi, S. Thomas, D. Dimitriadis, X. Cui, B. Ramabhadran, M. Picheny, L.-L. Lim, B. Roomi, and P. Hall, "English conversational telephone speech recognition by humans and machines," in *Conference of the International Speech Communication Association (INTERSPEECH)*, 2017, pp. 132–136.
- [23] Y. Tu, M.-W. Mak, and J.-T. Chien, "Variational domain adversarial learning for speaker verification," in *Conference of the International Speech Communication Association (INTERSPEECH)*, 2019, pp. 4315–4319.
- [24] S. Ioffe and C. Szegedy, "Batch normalization: Accelerating deep network training by reducing internal covariate shift," in *International Conference on Machine Learning (ICML)*, 2015, pp. 448–456.
- [25] C. Busso, M. Bulut, C.-C. Lee, A. Kazemzadeh, E. Mower, S. Kim, J. N. Chang, S. Lee, and S. S. Narayanan, "IEMOCAP: Interactive emotional dyadic motion capture database," *Language Resources and Evaluation*, vol. 42, no. 4, p. 335, 2008.
- [26] F. Ringeval, A. Sonderegger, J. Sauer, and D. Lalanne, "Introducing the RECOLA multimodal corpus of remote collaborative and affective interactions," in *IEEE International Conference and Workshops on Automatic Face and Gesture Recognition (FG)*, 2013, pp. 1–8.
- [27] D. P. Kingma and J. Ba, "Adam: A method for stochastic optimization," in *International Conference on Learning Representations (ICLR)*, 2015.
- [28] K. Sridhar, S. Parthasarathy, and C. Busso, "Role of regularization in the prediction of valence from speech," in *Conference of the International Speech Communication Association (INTERSPEECH)*, 2018, pp. 941–945.
- [29] B. Desplanques and K. Demuynck, "Cross-lingual speech emotion recognition through factor analysis," in *Conference of the International Speech Communication Association (INTERSPEECH)*, 2018, pp. 3648–3652.
- [30] S. Wu, T. H. Falk, and W.-Y. Chan, "Automatic speech emotion recognition using modulation spectral features," *Speech Communication*, vol. 53, no. 5, pp. 768–785, 2011.
- [31] S. Parthasarathy, R. Cowie, and C. Busso, "Using agreement on direction of change to build rank-based emotion classifiers," *IEEE/ACM Transactions on Audio, Speech, and Language Processing*, vol. 24, no. 11, pp. 2108–2121, 2016.
- [32] G. Trigeorgis, F. Ringeval, R. Brueckner, E. Marchi, M. A. Nicolaou, B. Schuller, and S. Zafeiriou, "Adeieu features? End-to-end speech emotion recognition using a deep convolutional recurrent network," in *IEEE international conference on acoustics, speech and signal processing (ICASSP)*, 2016, pp. 5200–5204.