

Unsupervised Personalization of an Emotion Recognition System: The Unique Properties of the Externalization of Valence in Speech

Kusha Sridhar, *Student Member, IEEE*, and Carlos Busso, *Senior member, IEEE*

Abstract—The prediction of valence from speech is an important, but challenging problem. The expression of valence in speech has speaker-dependent cues, which contribute to performances that are often significantly lower than the prediction of other emotional attributes such as arousal and dominance. A practical approach to improve valence prediction from speech is to adapt the models to the target speakers in the test set. Adapting a *speech emotion recognition* (SER) system to a particular speaker is a hard problem, especially with *deep neural networks* (DNNs), since it requires optimizing millions of parameters. This study proposes an unsupervised approach to address this problem by searching for speakers in the train set with similar acoustic patterns as the speaker in the test set. Speech samples from the selected speakers are used to create the adaptation set. This approach leverages transfer learning using pre-trained models, which are adapted with these speech samples. We propose three alternative adaptation strategies: unique speaker, oversampling and weighting approaches. These methods differ on the use of the adaptation set in the personalization of the valence models. The results demonstrate that a valence prediction model can be efficiently personalized with these unsupervised approaches, leading to relative improvements as high as 13.52%.

Index Terms—Speech emotion recognition, adaptation, transfer learning, emotional dimensions, valence.

◆

1 INTRODUCTION

THE area of *speech emotion recognition* (SER) is an important research problem due to its key potential in fields such as *human-computer interactions* (HCIs), healthcare [1], [2] and behavioral studies [3], [4]. Despite remarkable advances in emotion recognition, detecting emotions from speech is still a challenging task. The usual formulation to describe emotions is with categorical descriptors such as happiness, sadness, anger and neutral. However, this approach may not capture the intra and inter class variability across distinct emotional classes (i.e., variability across sentences with the same emotional class labels and variability across sentences with different emotional class labels). An alternative representation is the use of emotional attributes, as suggested by the core affect theory [5]. The most common attributes are arousal (calm versus active), valence (unpleasant versus pleasant) and dominance (weak versus strong). Because of their direct application in many areas, it is very important to build accurate models, which can reliably predict these emotional attributes. The estimation of emotional attributes is often posed as

a regression problem, where the goal is to predict the scores associated with these attributes. In particular, the emotional attribute valence is key to understand many behavioral disorders [6], [7] such as *post-traumatic stress disorder* (PTSD), depression, schizophrenia and anxiety. Although different approaches have been proposed to improve SER systems, the prediction of valence using acoustic features is often less accurate than other emotional attributes such as arousal or dominance. The gap in performance is significant even with different methods specially implemented to correct this problem, such as using features from other modalities [8], [9], [10], [11], modeling contextual information [12] or regularizing *deep neural network* (DNNs) under a *multitask learning* (MTL) framework [13], [14]. It is important to explore why predicting valence from speech is so difficult, and use these findings and insights to improve SER systems.

In our previous work, we studied the prediction of valence from speech [15], focusing the analysis on the role of regularization in DNNs. In particular, we explored the role of dropout as a form of regularization and analyzed its effect on the prediction of valence. Our analysis showed that a higher dropout rate (i.e., higher regularization) led to improvements in valence predictions. The

-
- K. Sridhar and C. Busso are with the Erik Jonsson School of Engineering and Computer Science, The University of Texas at Dallas, (E-mail: Kusha.Sridhar@utdallas.edu, busso@utdallas.edu)

optimal dropout rate for valence was higher than the optimal dropout rates for arousal and dominance across different configurations of the DNNs. A hypothesis from this study was that a heavily regularized network learns features that are more consistent across speakers, placing less emphasis on speaker-dependent emotional cues. We also conducted controlled speaker-dependent experiments to evaluate this hypothesis, where data from the same speakers were included in the train and test partitions. For valence, we observed relative gains in *concordance correlation coefficient* (CCC) up to 30% between speaker-dependent and speaker-independent experiments. The corresponding relative improvements observed for arousal and dominance were less than 4%. These results showed that valence emotional cues include more speaker-dependent traits, explaining why heavily regularizing a DNN helps to learn more general emotional cues across speakers [15]. Building on these results, we propose an unsupervised personalization approach that is very useful for valence prediction.

This paper explores the speaker-dependent nature of emotional cues in the expression of valence. We hypothesize that a regression model trained to detect valence from speech can be adapted to a target speaker. The goal is to leverage the information from the emotional cues of speakers in the train set to fine-tune a regression model already trained to perform well on the prediction of valence. Our approach identifies speakers in the train set that are closer in the acoustic space to the speakers in the test set. Data from these selected speakers are used to create an adaptation set to personalize the SER models toward the test speakers. We achieve the adaptation by using three alternative methods: unique speaker, oversampling and weighting approaches. The unique speaker approach randomly selects samples from the data obtained from the selected speakers in the train set without replacement, regardless of how many times these speakers are selected (i.e., a speaker in the train set may be found to be closer to more than one speaker in the test set). The oversampling approach draws data from the selected speakers as a function of the number of times that a given speaker is selected. For instance, if a speaker in the train set is found to be closer to two speakers in the test set, the selected sentences from that training speaker is counted twice. This approach repeats the data from this speaker during the adaptation phase, so the model sees the same speech samples in multiple batches in a single epoch. The weighting approach uses weights, where

samples from the selected speakers in the train set are weighted more. This approach adds weights on the cost function during the training process, building the models from scratch. We demonstrate the idea of personalization under two scenarios: 1) separate SER models, where each of them is personalized to a single test speaker (i.e., individual adaptation models), and 2) a single SER model personalized to a pool of 50 target speakers (i.e., global adaptation model).

Using the proposed model adaptation strategies leads to relative improvements in CCC as high as 13.52% in the prediction of valence. This result indicates the need for a personalization method to improve the prediction of valence, highlighting the benefits of our proposed approach. The contributions of our study are:

- We leverage the finding that the expression of valence in acoustic features is more speaker-dependent than arousal and dominance, raising awareness on the need for special considerations in its detection.
- We successfully personalize a SER system using unsupervised adaptation strategies by exploiting the speaker-dependent traits.
- We propose three alternative adaptation strategies to personalize a SER system, obtaining important relative performance improvements in the prediction of valence.

One of the key strengths of this study is that we find similar speaker in the emotional feature space alone. By exploiting similarities in the emotional feature space, we have suppressed the speaker trait or text dependencies. Our approach provides a much more powerful way of comparing emotional similarities than traditional methods used for speaker identification. Likewise, our personalization approach avoids or minimizes “concept drift.” SER is a challenging problem, where the prediction models can become more volatile with the addition of more data over time. Therefore, models built for analyzing such data quickly become obsolete. With our personalization study, we can minimize the impact of concept drift by developing personalized SER models that are tailored to target speakers. We can periodically update the models to target speakers or even weight the data based on their historical significance to develop better personalized models.

The paper is organized as follows. Section 2 discusses relevant studies on the prediction of valence from speech. It also describes the adaptation and personalization approach proposed for improving SER systems. Section 3 presents the database used

in this study. Section 4 describes the analysis on the role of regularization in the prediction of valence from speech, summarizing the study presented in our preliminary work [15]. Section 5 presents the proposed formulation to personalize a SER system, building on the insights learned from the analysis in Section 4. Section 6 presents the results obtained by using our proposed approaches to personalize a SER system. We primarily present the results on the MSP-Podcast corpus, but we evaluate the generalization of our proposed approach with two other databases. The paper concludes with Section 7, which summarizes our key findings, providing future directions for this study.

2 RELATED WORK

2.1 Improving the Prediction of Valence

While valence is a key dimension to understand complex human behaviors, its predictions using speech features are often lower than the predictions of other emotional attributes such as arousal or dominance [16]. Therefore, several speech studies have focused on understanding and improving valence prediction. Busso and Rahman [17] studied acoustic properties of emotional cues that describe valence. They built separate *support vector regression* (SVR) models trained with different groups of acoustic features: energy, fundamental frequency, voice quality, spectral, *Mel-frequency cepstral coefficients* (MFCCs) and RASTA features. They also built binary classifiers to distinguish between two groups of sentences characterized by similar arousal but different valence. The study showed that spectral and fundamental frequency features are the most discriminative for valence. Koolagudi and Rao [18] claimed that MFCCs were effective to classify emotion along the valence dimension (i.e., spectral features). Cook et al. [19], [20] explored the structure of the fundamental frequency (F0), extracting dominant pitches in the detection of valence from speech. Deshpande et al. [21] proposed a reduced feature set consisting of the autocorrelation of pitch contour, *root mean square* (RMS) energy and a 10-dimensional *time domain difference* (TDD) vector. The TDD vector corresponds to successive differences in the speech signal. The feature set collectively led to better results than MFCCs or OpenSmile features [22]. Tur-sunov et al. [23] used acoustic descriptors associated with timbre perception to classify discrete emotions, and emotions along the valence dimension. Tahon et al. [24] showed that voice quality features were also useful in the detection of valence.

Other studies have explored modeling strategies to improve the prediction of valence. Lee et al. [25] used dynamic Bayesian networks to capture time dependencies and mutual influence of interlocutors during dyadic interactions. Contextual information was found to be particularly useful in the prediction of valence, leading to relative improvements higher than the one observed for arousal. Another alternative approach to improve valence was by regularizing a DNN. For example, Parthasarathy and Busso [13] showed that jointly predicting valence, arousal and dominance under a *multitask learning* (MTL) framework helps to improve its prediction. The MTL framework acts as regularization in DNNs. Other approaches using MTL have shown similar findings. Since the performance of lexical models often outperforms acoustic models in predicting valence [9], Lakomkin et al. [26] suggested the use of the output of an *automatic speech recognition* (ASR) as the input of a character-based DNN.

2.2 Model Adaptation in SER Tasks

Unlike other speech tasks such as *automatic speech recognition* (ASR) that rely on abundant data, databases used in SER are often small. Therefore, many researchers have explored the use of model adaptation techniques to generalize the models beyond the training conditions. Most of the adaptation techniques aim to attenuate sources of variability including channel, language and speaker mismatches. Early studies demonstrated the effectiveness of these techniques with algorithms based on *support vector machine* (SVM) [27]. Abdelwahab and Busso [28] demonstrated the importance of data selection strategy for domain adaptation. They illustrated that, incrementally adapting emotion classification models using active learning to select samples from the target domain can improve their performance. They used a conservative approach where only the correctly classified samples were used to adapt the model, leaving out the incorrect ones to avoid large changes in the hyperplane between the classes.

Recent efforts in model adaptation have mainly focused on DNNs, where important advances have been made in the area of transferring knowledge between domains [29]. DNNs with their deep architectures can learn useful representations by compactly representing functions. Deng et al. [30] used sparse autoencoders to learn feature representations in the source domain that are more consistent with the target domain. This goal was achieved by simultaneously minimizing the reconstruction error in both domains. Deng et al. [31] proposed the use

of unlabeled data under a deep autoencoder framework to reduce the mismatch between train and test conditions. They also simultaneously learned common traits from labeled and unlabeled data.

Instead of the traditional method of pre-training and fine-tuning for model adaptation, Gideon et al. [32] used progressive networks to enhance a SER system. They trained the model on new tasks by freezing the layers related to previously learned tasks and used their intermediate representations as inputs to new parallel layers. This study also used paralinguistic information from gender and speaker identity to achieve improvements. Similarly, other variants of adaptation techniques use *kernel mean matching* (KMM) [33], *Nonnegative matrix factorization* [34], *domain adaptive least-squares regression* [35], and PCANet [36]. These methods lead to improvements on emotion recognition tasks, by using hybrid frameworks involving unsupervised followed by supervised learning. Our proposed approach is different from these studies, since we aim to explicitly exploit similarities between speakers in the train and test sets, as measured in the feature space. This approach leads to powerful adaptation methods that are particularly useful to predict valence.

2.3 Speech Emotion Personalization

This study focuses on adapting or personalizing a SER system to a target set of speakers. Busso and Rahman [37] demonstrated the idea of personalization using an unsupervised feature normalization scheme. They used the *iterative feature normalization* (IFN) method [38] to reduce speaker variability, while preserving the discriminative information of the features across emotional classes. The IFN algorithm has two steps. First, it detects neutral sentences which are used to estimate the normalization parameters. Then, the data is normalized with these parameters. Since the detection of neutral speech is not perfect, this process is iteratively repeated leading to important improvements. Busso and Rahman [37] implemented the IFN scheme as a front end of a SER system designed to recognize emotion from a target speaker, observing better accuracy.

3 RESOURCES

3.1 Emotional Corpora

3.1.1 The MSP-Podcast Corpus

The study relies on the MSP-Podcast corpus [39], which provides a diverse collection of spontaneous speech segments that are rich in emotional content. The speech segments are obtained from podcasts

taken from various audio-sharing websites, using the retrieval-based approach proposed by Mari-ooryad et al. [40]. The content of the podcasts is diverse, including discussions on sport, politics, entertainment, games, social problems and healthcare. The podcasts are segmented into speaking turns between 2.75s and 11s duration. These segments are automatically processed to discard segments with music, overlapped speech, and noisy recordings. Since most of the segments are expected to be neutral, we retrieve candidate segments to be included in the database by leveraging a diversified set of SER algorithms to detect emotions. The selected speech segments are annotated on *Amazon Mechanical Turk* (AMT) using a crowdsourcing protocol similar to the one introduced by Burma-nia et al. [41]. This crowdsourcing protocol stops the annotators in real-time if their performance is evaluated as poor. The raters annotate each speaking turn for its arousal, valence and dominance content using *self-assessment manikins* (SAMs) on a seven Likert-type scale. The ground truth labels for each speaking turn is the average across the scores provided by the annotators. Although we do not use categorical annotations in this study, the corpus also includes annotations of primary and secondary emotions. The primary emotion corresponds to the dominant emotional class. The secondary emotion corresponds to all the emotional classes that can be perceived in the speech segments.

The collection of the MSP-Podcast corpus is an ongoing effort. This study uses version 1.6 of the MSP-Podcast corpus, which consists of 50,362 speech segments (83h29m) annotated with emotional classes. From this set, 42,567 segments have been manually assigned to 1,078 speakers. The speaker identity for the rest of the corpus has not been assigned yet. Figure 1 illustrates the partition of the dataset used in this study. The test set has 10,124 speech segments from 50 speakers, and the development set has 5,958 speech segments from 40 speakers. Each speaker in the test and development sets has a minimum of five minutes of data. The rest of the corpus is included in the train set, which consists of a total of 34,280 speech segments. The data partition aims to create speaker-independent partitions between sets. Lotfian and Busso [39] provide more details on this corpus.

As shown in Figure 1, we further split the test set into two partitions for this study: *test-A* and *test-B* sets. The *test-A* set includes 200s of recording for each of the 50 speakers in the test set. The *test-B* set includes the rest of the recordings in the test set.

Each test speaker has at least 300s (5 mins) of data. After removing 200s from each speaker to form the *test-A* set, the *test-B* set is left with at least 100s of data for each speaker. The average duration per speaker in the *test-B* set is 1005.96s.

3.1.2 The IEMOCAP and MSP-IMPROV Corpora

Besides the MSP-Podcast corpus, we use two other databases for our experimental evaluations. The first database is the USC-IEMOCAP corpus [42], which is an audiovisual corpus and contains dyadic interactions from 10 actors in improvised scenarios. This study only uses the audio. The database contains 10,039 speaking turns, which are annotated with emotional labels for arousal, valence and dominance by at least two raters using a five-Likert scale. We also use the MSP-IMPROV corpus [43], which is a multimodal emotional database that contains interactions between pairs of actors engaged in improvised scenarios. In addition to the conversations during improvised scenarios, the dataset also contains the interactions between the actors during the breaks, resulting in more naturalistic data. The corpus uses a novel elicitation scheme, where two actors in an improvised scenario leads one of them to utter target sentences. For each of the target sentences, four emotional scenarios were created to contextualize the sentence to elicit happy, angry, sad and neutral reactions, respectively. This corpus consists of 8,438 turns of emotional sentences recorded from 12 actors (over 9 hours). The sessions were manually segmented into speaking turns, which were annotated with emotional labels using perceptual evaluations. Each turn was annotated for arousal, valence and dominance by five or more raters using a five-Likert scale. In both databases, the consensus emotional attribute label assigned to each utterance is the average across the scores provided by the annotators, which is linearly mapped between -3 and 3.

3.2 Acoustic Features

This study uses the feature set proposed for the *computational paralinguistics challenge* (ComParE) in Interspeech 2013 [44]. The features are extracted by estimating several *low-level descriptors* (LLDs) such as energy, fundamental frequency and MFCCs. For each speech segment, statistics such as mean, standard deviation, range and regression coefficients are estimated for each LLD, creating *high-level descriptors* (HLDs). With this approach, the feature vector is fixed regardless of the duration of the sentence. The ComParE set creates a 6,373 dimensional feature vector for each sentence.

MSP-Podcast Corpus

Train set	Development set	Test set				
34,280 sentences At least 1,078 speakers	5,958 sentences 40 speakers	10,124 sentences 50 speakers				
		<table border="1"> <thead> <tr> <th>Test-A</th> <th>Test-B</th> </tr> </thead> <tbody> <tr> <td>200s per speaker 50 speakers</td> <td>Rest of test set 50 speakers</td> </tr> </tbody> </table>	Test-A	Test-B	200s per speaker 50 speakers	Rest of test set 50 speakers
Test-A	Test-B					
200s per speaker 50 speakers	Rest of test set 50 speakers					

Fig. 1. Partitions of the MSP-Podcast corpus used in this study for the train, development and test sets. The test set is further split into the *test-A* and *test-B* sets.

4 ROLE OF REGULARIZATION

Our proposed personalization method for valence builds upon the findings reported in our preliminary study [15]. This section summarizes the main findings on the role of dropout rate as a form of regularization in DNNs and its impact on SER. The study in Sridhar et al. [15] was conducted on an early version of this corpus. We update the analysis with the release of the corpus used for this study (release 1.6 of the MSP-Podcast corpus).

4.1 Optimal Dropout Rate for Best Performance

Our previous study focused on the role of dropout as a form of regularization in improving the prediction of valence [15]. When dropout is used in DNNs, random portions of the network are shutdown at every iteration, training a smaller network on each epoch. This approach helps in learning feature weights in random conjunction of neurons, preventing developing co-dependencies with neighboring nodes. This regularization approach leads to better generalization. We explore the role of regularization in the prediction of valence by changing the dropout rate p . The goal of this analysis is to understand the optimal value p that leads to the best performance for different network configurations (i.e., different number of layers, different number of nodes per layer). We train the models for 1,000 epochs, with an early stopping criterion based on the development loss. The loss function is based on CCC, which has led to better performance than *mean squared error* (MSE) [16]. We train separate regression models by changing the dropout rate $p \in \{0.0, 0.1, \dots, 0.9\}$, recording the optimal dropout rate leading to the best performance on the development set. We evaluate two networks with three and seven layers, implemented with 256, 512 and 1,024 nodes per layer. Note that deeper models have dropout in more layers, so more nodes are being turn off. Therefore, comparisons of dropout rates on networks with different layers are less meaningful than comparisons of the optimal dropout rate across emotional

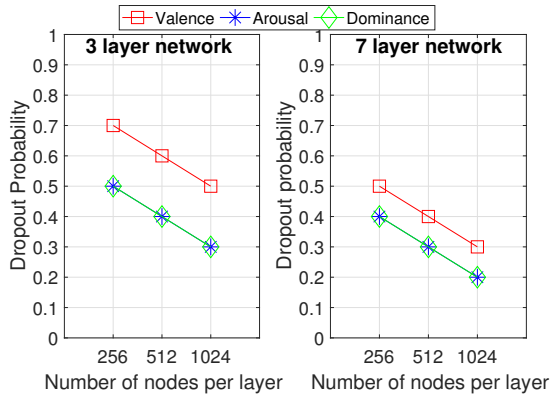


Fig. 2. Optimal dropout rate observed in the development set as a function of the number of nodes per layer in a DNN. The DNN is implemented with either three or seven layers.

attributes. Figure 2 illustrates the results, showing the optimal dropout rate observed in the development set. The optimal dropout rates that give the best performance are higher for valence than arousal and dominance. Figure 2 shows that the gap between the optimal dropout rates for valence and arousal/dominance stays consistent across conditions. Interestingly, the optimal dropout rates for arousal and dominance are exactly the same across different DNN configurations, whereas it is different for valence. The results show that the need for higher regularization for valence is consistent across variations in the architectures of the DNNs.

4.2 Speaker-Dependent versus Speaker-Independent Models

Section 4.1 demonstrated that a DNN needs to be heavily regularized to give good predictions for valence. We hypothesize that this finding can be explained by the speaker-dependent nature of speech cues in the expression of valence (i.e., we use different acoustic cues to express valence). A DNN with higher regularization, learns more generic trends present across speakers, leading to better generalization. To validate our hypothesis, we conduct a controlled emotion detection evaluation, where we train DNNs with either speaker-dependent or speaker-independent partitions. SER should be performed with speaker-independent partitions, where the data in the train and test partitions are from disjoint set of speakers. A model trained with data from speakers in the test set has an unfair advantage over a system evaluated with data from new speakers, resulting in overestimated performance. Our goal is to quantify the benefits of using speaker-dependent partitions.

We build DNNs with four layers, implemented with 256, 512 or 1,024 nodes. The speaker-

TABLE 1

Comparison of CCC values between speaker-independent and speaker-dependent conditions. The DNN is trained with four layers. The column 'Gain' shows the relative improvement by training with partial data from the target speakers (*test-A* set).

Attributes	Nodes	Speaker Independent	Speaker Dependent	Gain (%)
		<i>Test-B</i> set	<i>Test-B</i> set	
Valence	256	0.3076	0.3373	9.65
	512	0.3083	0.3670	19.03
	1,024	0.2997	0.3538	18.05
Arousal	256	0.7153	0.7216	0.88
	512	0.7164	0.7331	2.33
	1,024	0.7104	0.7258	2.16
Dominance	256	0.6300	0.6379	1.25
	512	0.6374	0.6565	2.99
	1,024	0.6253	0.6352	1.58

dependent model is built by adding the *test-A* set in the train set (Fig. 1). This approach creates a train set with partial knowledge about the speakers in the test set. In contrast, the speaker-independent model uses only the train set. To have a fair comparison, both models are evaluated on the *test-B* set with speech samples that are not used to either train or optimize the parameters of the systems. Table 1 shows the CCC values of the models for speaker-independent and speaker-dependent conditions. The last column calculates the relative improvements achieved under the speaker-dependent condition. We observe up to 19.03% gain in performance for valence. The performances for the arousal and dominance models also increase, but the relative improvements are less than 3%. These results are not due to the extra samples added in the train set. The *test-A* set only represents a 4.4% of the size of the train set. In fact, evaluations where we randomly remove data from the train set to maintain a consistent size led to similar performance than the speaker dependent condition reported in Table 1. The fact that the performances increase using speaker dependent sets is expected. What is unexpected is that the relative gain is significantly higher for valence than for arousal and dominance. These results clearly show that learning emotional traits from the target speakers in the test set has clear benefits for valence, validating our hypothesis that the expression of valence in speech has speaker-dependent traits.

5 PROPOSED PERSONALIZATION METHOD

5.1 Motivation

The findings in Section 4 suggest that leveraging data from speakers in the train set that are *closer* to our target speakers in the test set should benefit our SER models. This is the premise of our proposed

approach. We aim to improve the prediction of valence, bridging the gap in performance between speaker-dependent and speaker-independent conditions reported in Table 1. Data sampled from the selected speakers' recordings are used to create an adaptation set, as illustrated in Figure 3(a). Once the closest speakers are identified, we can either adapt the models or assign more weights to samples from this adaptation set. This section describes our unsupervised personalization approach to improve the prediction of valence. Unlike the speaker-dependent settings used in Section 4.1 (and Sec. 6.8), the analysis and experiments in this study operate with speaker-independent partitions for train, development and test sets. The assumption in our formulation is that we have the speaker identity associated with each sentence in the test set.

5.2 Estimation of Similarity Between Speakers

A key step in our approach is to identify speakers in the train set that are *closer* to the speakers in the test set. Ideally, we would like to identify speakers who externalize valence cues in speech in a similar way. This aim is difficult with no clear solutions. We simplify our formulation by searching for similarities between speakers in the space of emotional speech features. By exploiting similarities on the emotional feature space, we expect to focus more on emotional patterns than on speaker traits, which would be the focus of speaker embeddings created with methods such as the i-vector [45] or the x-vector [46]. Our approach relies on *principal component analysis* (PCA) to reduce the dimension of the space, followed by fitting a *Gaussian mixture model* (GMM) to the resulting reduced feature space.

We aim to quantify the similarity in the feature space between the speaker i in the train set, and the speaker j in the test set, $d(i, j)$. The first step is to reduce the feature space, since we consider a high dimensional feature vector (6,373D – Sec. 3.2). Reducing the feature space creates a more compact feature representation, where the similarity between speakers can be more efficiently computed. We implement this step with PCA, which is a popular unsupervised dimensionality reduction technique. First, we estimate the zero-mean vector $\mathbf{y}_s = \mathbf{f}_s - \bar{\mathbf{f}}$, where \mathbf{f}_s is the feature vector of sentence s , and $\bar{\mathbf{f}}$ is the mean feature vector. Then, we concatenate these M vectors, creating matrix F (Eq. 1). From this matrix, we estimate the sample covariance matrix Q using Equation 2. Then, we compute the eigenvectors of Q , selecting the ones with the highest

eigenvalues, which are considered as the *principal components* (PCs).

$$F = [\mathbf{y}_1, \mathbf{y}_2, \dots, \mathbf{y}_M] \quad (1)$$

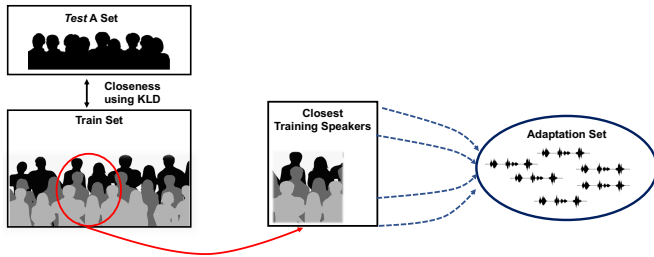
$$Q = \frac{1}{M-1} F F^T \quad (2)$$

The PCA-based feature reduction is implemented for each speaker in the test set, creating speaker-dependent transformations. We use the 10 most important dimensions, which explain in average 57.9% of the variance in the feature space. The speech sentences from speaker i (train set) are projected into the PCA space associated with speaker j (test set). The speech sentences from speaker j are also projected in that space. After the PCA projections, we fit two separate GMMs on the reduced feature space, one for the sentences of speaker i (p_i^{train}), and another for the sentences of speaker j (q_j^{test}). We empirically set 10 mixtures for the GMMs (there is no direct connection between the number of mixture components used for the GMMs and the dimensionality of the PCA space). Finally, we estimate the similarity between the GMMs using the *Kullback Liebler Divergence* (KLD). We employ a numerical approximation using multiple Monte Carlo simulations. We sample a set of data points from each pair of GMMs to compute the KLD score. We repeat this procedure multiple times, estimating the average results.

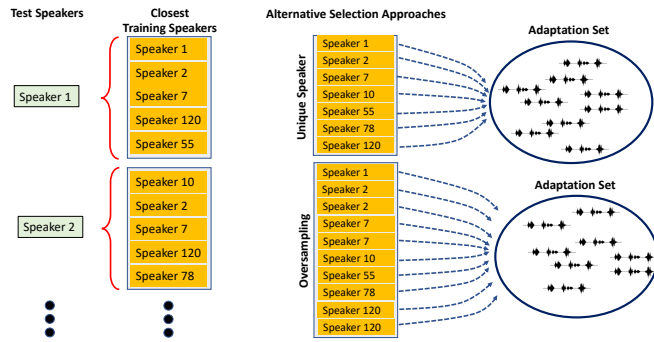
For a given speaker j in the test set, we estimate $d(i, j) = KLD(p_i^{train}, q_j^{test})$ for all the speakers in the train set, sorting their scores in increasing order. The closest speakers in the train set are the top speakers in this ranked list. This approach is repeated for each of the 50 speakers in the test set (i.e., we have 50 different PCA projections). While this step can be implemented using all the data from the test set, we use the *test-A* set to have the same amount of data for each speaker (i.e., 200s – Fig. 1). Figure 3(a) illustrates the process to form the adaptation set by finding the closest set of training speakers to a target speaker. Notice that the adaptation set is a subset of the train set, for which we have labels.

5.3 Personalization Approach

After selecting the speakers in the train set that are closest to each of the test speakers, our next step is to leverage data from these train set speakers to either personalize or adapt the emotion prediction models for valence. We propose and evaluate three alternative methods, referred to as *unique speaker*, *oversampling*, and *weighting* approaches. The first two approaches rely on adapting a model. We built a regression model using the train set to predict



(a) Selection of closest speakers to create adaptation set



(b) Criteria for selecting samples for the adaptation set

Fig. 3. Illustration of proposed personalization approach. We identify the closest set of speakers in the train set to each of the target speaker in the test set. Sentences from these speakers are randomly sampled with different selection criteria.

valence. The weights of the pre-trained models are frozen with the exception of the last layer, which are fine-tuned using the adaptation data. The last approach requires training the network from scratch.

Unique speaker approach: Each speaker in the test set generates a list with its N closest speakers in the train set. Since some speakers in the train set may be close to more than one speaker, the total number of selected speakers after combining the lists from the 50 speakers in the test set is less or equal to $N \times 50$. The unique speaker approach considers all the data from these speakers in the train set. Each speech segment can be considered only once in the adaptation set. The order in which we consider the target speakers does not affect the final selection of the data, since the algorithm selects for each speaker in the test set her/his closest speakers from the entire set of speakers in the train set. Figure 3(b) illustrates the process for the case when we have only 2 speakers in the test set. In the example, speakers 2, 7 and 120 in the train set are found to be close to both test speakers, hence, we consider them only once when forming the adaptation set. We implement a balance sampling criterion that aims to select approximately the same amount of data for each speakers. For example, for an adaptation set of 200s, if we have 7 unique speakers selected as

the closest speakers from the test speakers, as in the example, we would randomly select approximately 28.6s for each of these speakers. We adopt this approach in an attempt to diversify and balance the speech samples selected from all the speakers in the unique speakers set. This approach uses the pre-trained models built with all the train set, personalizing the models with the adaptation set.

Oversampling approach: A speaker in the train set may be in the list of the closest speakers for more than one speaker in the test set. The oversampling approach assumes that these samples are more relevant during the adaptation process. If a speaker is selected C times (i.e., the speaker is one of the closest speakers for C speakers in the test set), the oversampling method will create C copies of his/her sentences before randomly drawing the samples. This process is illustrated in Figure 3(b) where speakers 2, 7 and 120 in the train set are copied twice. Therefore, more samples from these speakers will appear in the adaptation set. We form the adaptation set with a balance approach, choosing approximately the same amount of data from the selected speakers. In the example from Figure 3(b), we would select 20 second for each of the 10 sets for an adaptation set of 200s. Sentences can even be repeated on the adaptation set. This approach also fine-tunes the pre-trained model using speech samples from the oversampled adaptation set.

Weighting approach: The third approach to personalize a model is by increasing the weights in the loss function during the training process for speech samples in the adaptation set (i.e., same set used in the unique speaker approach). Unlike the previous two approaches, which adapt a pre-trained system, this approach trains the regression model from scratch. As described in Equation 3, we use $\mathcal{L} = (1 - CCC)$ as the loss function to train our models. For the weighting approach, this cost is assigned to a sample in the train set, but not on the adaptation set. For samples in the adaptation set, we multiply the cost \mathcal{L} by a factor $\lambda > 1$. Therefore, an error of a sample from the adaptation set is λ times more costly than an error made on other samples from the train set not included in the adaptation set. We experiment with weighting ratios of 1:2, 1:3, 1:4 and 1:5, where higher weights are assigned to samples in the adaptation set. This approach uses all the train set, increasing the importance of correctly predicting samples in the adaptation set.

The proposed unsupervised adaptation schemes can be jointly applied to all the speakers in the test set, creating a single model. We refer to this

approach as *global adaptation* (GA) model. Alternatively, the approaches can be individually implemented for each speaker, creating as many models as speakers in the test set. This implementation only works for the unique speaker and weighting approaches. The oversampling approach does not apply in this case, since each speech segment in the adaptation set is drawn only once (i.e., we consider one test speaker at a time). We refer to this approach as the *individual adaptation* (IA) model. We evaluate both implementations in Section 6 using adaptation sets of different sizes.

6 EXPERIMENTAL RESULTS

The prediction of valence is formulated as a regression problem implemented with DNNs with four dense layers and 512 nodes per layer. This setting achieved the best performance for valence on Table 1. We use *rectified linear unit* (ReLU) at the hidden layers and linear activations for the output layer. The DNNs are trained with batch normalization for the hidden layers. We use a dropout rate of $p = 0.7$ at the hidden layers. The selection of this rate follows the findings in Section 4.1, which demonstrates that a higher value of dropout is important for improving the detection of valence. We pre-train the models for 200 epochs with an early stopping criterion based on the performance on the development set. The best model is used to evaluate the results. The DNNs are trained with *stochastic gradient descent* (SGD) with momentum of 0.9, and a learning rate of $r = 0.001$. For the unique speaker and oversampling approaches, the learning rate is reduced to $r_{adapt} = 0.0001$ while adapting the regression model. We adapt the models with these approaches for 100 extra epochs with an early stopping criterion based on the performance on the development set. For the weighting approach, we train the models from scratch for 200 epochs with early stopping criterion, maximizing the performance on the development set. The loss function (Eq. 3) relies on the *concordance correlation coefficient* (CCC), which is defined in Equation 4.

$$\mathcal{L} = (1 - CCC) \quad (3)$$

$$CCC = \frac{2\rho\sigma_x\sigma_y}{\sigma_x^2 + \sigma_y^2 + (\mu_x - \mu_y)^2} \quad (4)$$

The parameters μ_x and μ_y , and σ_x and σ_y are the means and standard deviations of the true labels (x) and the predicted labels (y), and ρ is the Pearson's correlation coefficient between them. CCC takes into account not only the correlation between the

true emotional labels and their estimates, but also the difference in their means. This metric takes care of the bias-variance tradeoff when comparing the true and predicted labels. CCC is also the evaluation metric in all our experimental evaluation.

The input to the DNNs is the 6,373D acoustic feature vector (Sec. 3.2). The features are normalized to have zero mean and unit standard deviation. This normalization is done using the mean and standard deviation values estimated over the training samples. After this normalization, we expect the features to be within a reasonable range. We remove outliers by clipping the features with values that deviate more than three standard deviations from their means (i.e., $\mu_{f_i} - 3\sigma_{f_i} \leq f_i \leq \mu_{f_i} + 3\sigma_{f_i}$). The output of the DNNs is the prediction score for the emotional attribute.

We use the speaker-independent and speaker-dependent models described in Section 4.2 as baselines, where the results are listed in Table 1. The speaker-independent model does not rely on any adaptation scheme to personalize the models to perform better on the test set. As described in Section 4.2, the speaker-dependent model is built by adding the *test-A* set to the train set, using partial information from the speakers. While this setting is not representative of the performance expected for the regression model when evaluated on speech from unknown speakers, it provides an upper bound performance to contextualize the improvements observed with our proposed personalization methods. For analysis purposes, we report the performance of the speaker-dependent models obtained with the addition of 50s, 100s, 150s, and 200s per speaker in the test set. These extra samples are obtained from the *test-A* set. We consistently evaluate all the models using the *test-B* set. This experimental setting creates fair comparisons, since these samples are not used to train any of the models.

6.1 Global Adaptation Model

We evaluate the performance of the system with the global adaptation model, where a single regression model is built. The three adaptation schemes are implemented by considering the 50 speakers in the test set. The adaptation set is obtained by identifying the closest speakers in the train set to each speaker in the test set (Sec. 5.2). We implement this approach by identifying the five closest speakers in the train set ($N = 5$). Section 6.3 shows the results with different number of speakers. We incrementally add more samples by randomly selecting 50s, 100s, 150s, 200s, and 300s from the selected speakers associated with

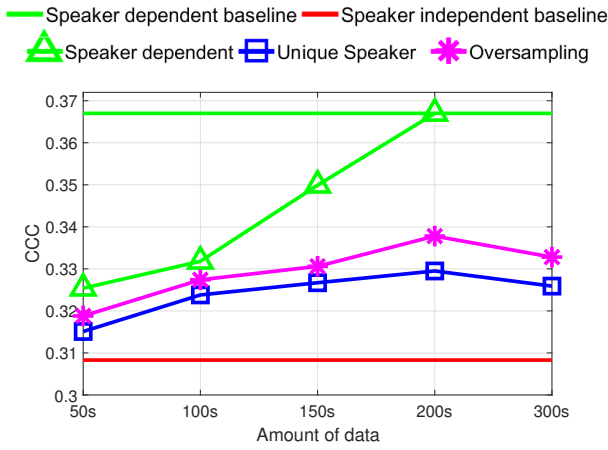


Fig. 4. Results on the *test-B* set for the global adaptation model for the unique speaker and oversampling methods. It shows the results for the speaker-dependent and speaker-independent baselines. It also shows the speaker-dependent model implemented with different size of data from the *test-A* set.

a given speaker in the test set (Sec. 5.3). This process is repeated for each of the speakers in the test set to observe the performance trend as a function of the size of the adaptation data.

First, we evaluate the unique speaker and oversampling approaches, which rely on model adaptation. Figure 4 shows the results. The two solid horizontal lines are the speaker-dependent (green) and speaker-independent (red) baselines. The green line (triangle) corresponds to the speaker-dependent model as we increase the amount of data from the *test-A* set. The performance for the unique speaker approach is shown in blue (square). We clearly observe an improvement over the speaker-independent baseline, which demonstrate that the adaptation scheme is effective even with a very small adaptation set (e.g., 50s). The pink (asterisk) line in Figure 4 shows the performance of the oversampling approach, which leads to better performance than the unique speaker approach. Both approaches use the same amount of adaptation data, but rely on different criteria to select the adaptation samples. Adding samples in multiple minibatches according to the oversampling strategy is beneficial to improve the prediction of valence. For both models, we observe consistent improvement in CCC as more data is added into the adaptation set from 50s to 200s. After this point, the performance seems to saturate, observing fluctuations. Interestingly, adaptation with 200s of data using the oversampling approach leads to better performance than the speaker-dependent baseline implemented with 50s and 100s of data from the *test-A* set.

Second, we evaluate the performance of the

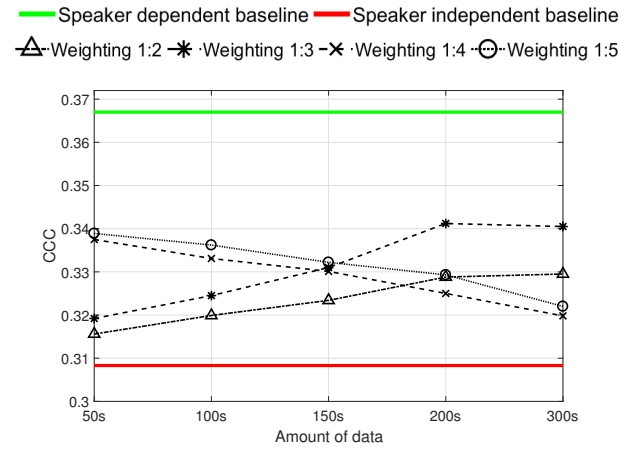
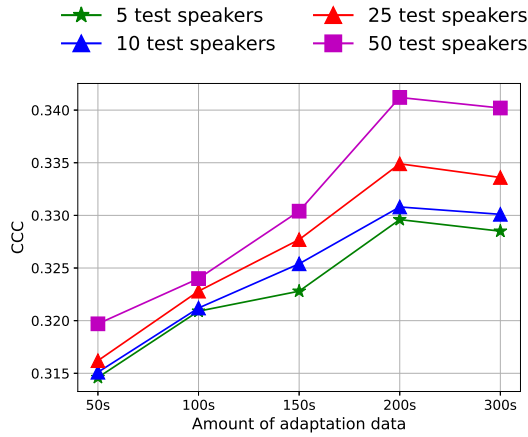


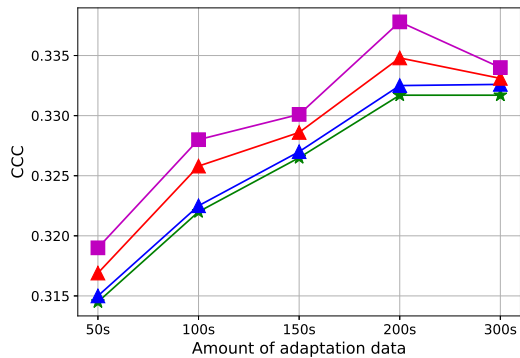
Fig. 5. Results on the *test-B* set for the global adaptation model for the weighting approach. The figure shows the results for the speaker-dependent and speaker-independent baselines.

weighting method, which trains the models from scratch, weighting more the samples in the adaptation set (Sec. 5.3). We evaluate the amount of data included in this selected set, including 50s, 100s, 150s, 200s per speaker in the test set. We also consider using all the data from the selected speakers. Only samples in this set are weighed more, implementing this approach with different ratios (1:2, 1:3, 1:4 and 1:5). Figure 5 shows the results. For weighting ratios 1:2 and 1:3, the performance gradually increases by adding more data in the selected set, peaking at 200s per speaker. However, the opposite trend is observed when the weighting ratios are either 1:4 or 1:5. Increasing the weights of speech samples in the adaptation set diminishes the information provided in the rest of the train data, leading to worse performance. The best performance is obtained with a weighting ratio of 1:3, when the selected set includes 200s per speaker. Figure 5 shows that this setting achieves similar performance than the best setting for the oversampling approach (Fig. 4).

The global adaptation models are trained with 50 target speakers. We also investigate the performance, when we reduce the number of target speakers. We implement the approach with 5, 10, 25 or 50 target speakers. This analysis only uses the two best models for the global adaptation approach, which are the weighting method with a 1:3 ratio, and the oversampling approach. Figure 6 shows that the performance trends are similar for these four conditions as we increase the amount of adaptation data. This result is expected, because we are still using the closest training speakers to each target speaker to perform the adaptation, which is



(a) Weighting approach



(b) Oversampling approach

Fig. 6. Performance of the global adaptation model with different number of target speakers as we increase the size of the adaptation set. The figure reports results using the weighting (1:3 ratio) and oversampling approaches on the *test-B* set.

shown to be beneficial for predicting valence. The best performance is achieved by using the 50 target speakers. The slight decrease in performance by using less target speakers may be explained by the reduced size of the adaptation set, noting that we add 50s, 100s, 150s, 200s, or 300s for each of the target speakers.

6.2 Individual Adaptation Model

This section presents the results of our approach implemented using the individual adaptation model. This approach builds one model for each of the speakers in the test set, creating the adaptation set with the samples from the speakers in the train set that are closer to this speaker (i.e., 50 separate models). For each model, we attempt to select equal duration of speech samples from each of the closest train speakers in the adaptation set to balance the amount of data used from each speaker. After adapting the models, the results are reported by concatenating the predicted vectors for each speaker in the test set. We estimate the CCC

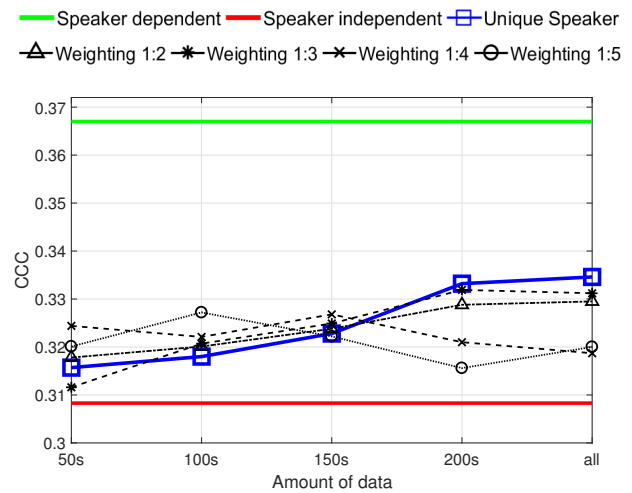


Fig. 7. Results on the *test-B* set for the individual adaptation model using the unique speaker and weighting approaches. The figure shows the results for the speaker-dependent and speaker-independent baselines.

values for the entire *test-B* set. The approach is implemented with the five closest speakers to each speaker in the test set. The performance of the approaches are reported by increasing the adaptation set. As explained in Section 5.3, we only evaluate the unique speaker and weighting approaches, since the oversampling approach cannot be implemented with a single speaker in the test set.

Figure 7 shows the CCC scores obtained for different sizes of the adaptation set. The results show improvements over the speaker-independent baseline performance. The performance gains are consistently higher when all the data from the closest speakers are used. The weighting approach also leads to better performance than the unique speaker approach. However, the results are worse than the CCC values of approaches implemented with the global adaptation model (Figs. 4 and 5). The decrease in performance in the IA model can be associated with the adaptation procedure. In the IA models, we adapt separate models for each target speaker. This procedure involves adapting 50 different models where their parameters and hyperplanes change based on a small adaptation set. This approach may be too aggressive, resulting in lower performance than adapting a single model, considering all the target speakers. We have seen similar observations in the area of active learning in SER tasks, where a more conservative adaptation strategy led to better results [28]. In the GA models, we are adapting a single model, where the shape and direction of the change of the hyperplane are smoother than in the case of the IA models achiev-

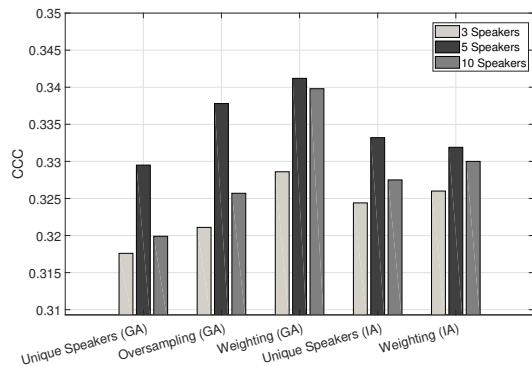


Fig. 8. Evaluation of the optimal number of closer speakers (N) from the train set for each of the speakers in the test set. The training data from the selected speakers is included in the adaptation set. The results corresponds to the CCC values obtained with different methods on the *test-B* set.

ing a better and stable performance.

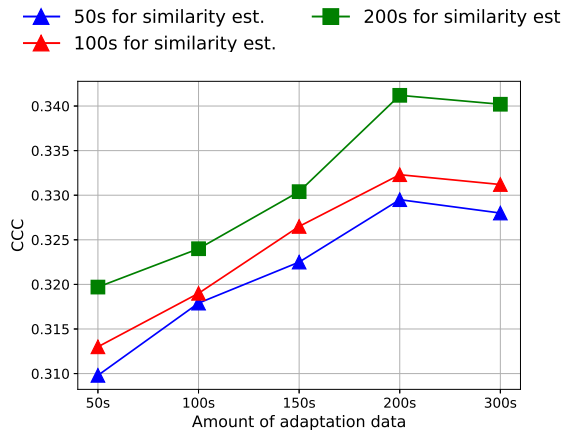
6.3 Number of Closest Speakers

This section evaluates the number of closest speakers (N) from the train set selected for each speaker in the test set. If this number is too small, the adaptation set will not have enough variability. If this number is too high, we will select speakers that are not very close to the target speakers. We implement the weighting approach with the ratio 1:3. Figure 8 shows the results on the *test-B* set for the models implemented with the global and individual adaptation models using the proposed adaptation schemes ($N \in \{3, 5, 10\}$). The results clearly show that $N = 5$ is a good balance, obtaining higher performance across conditions. We set $N = 5$ for the rest of the experimental evaluation.

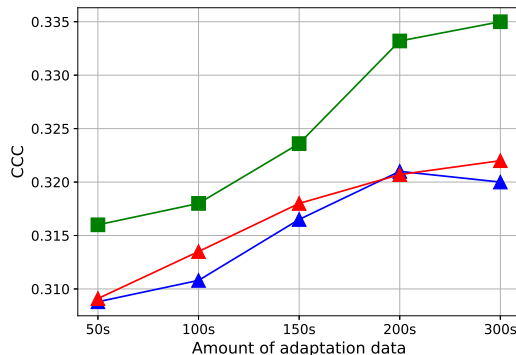
6.4 Size of Target Speaker Data for Speaker Similarity

Throughout the paper we have used 200s of data from each target speaker to perform speaker similarity estimation. This section evaluates the performance of the personalization method when we use less data from each target speaker. We consider the best performing method for the global (i.e., weighting method with a 1:3 ratio) and individual (i.e., unique speaker method) personalization formulations using 50s, 100s and 200s of data from each target speaker. We kept the *test-B* set consistent across all the experiments, without augmenting it with data from the *test-A* that we do not use.

Figure 9 shows the results obtained for the global (Fig. 9(a)) and individual (Fig. 9(b)) adaptation methods. The results show consistent patterns



(a) Global Adaptation Model – Weighting Approach



(b) Individual Adaptation Model – Unique Speaker Approach

Fig. 9. Adaptation results for different size of data from the target speakers used for estimating speaker similarity. The figure shows the results for the weighting method for the global adaptation experiments with a 1:3 ratio, and the unique speaker method for the individual adaptation experiments.

when we used either 50s, 100s or 200s of the target speakers. The adaptation methods lead to increase in performance as we add adaptation data draw from the closest speakers in the train set. The best performance for both global and individual adaptation methods is when we use 200s of data from the target speakers to estimate the speaker similarity estimation. Unfortunately, we cannot evaluate cases when we have more than 200s, given that we want to keep the *test-B* set consistent for all the evaluations. This evaluation suggests that a minimum of 200s of data from the target speakers is needed to successfully estimate the speaker similarity, and achieve the best performance.

6.5 Minimizing Loss on Adaptation Set

The results of the model adaptation presented in Sections 6.1, 6.2, 6.3 and 6.4 are obtained by minimizing the loss function on the development set, a practice that aims to increase the generalization

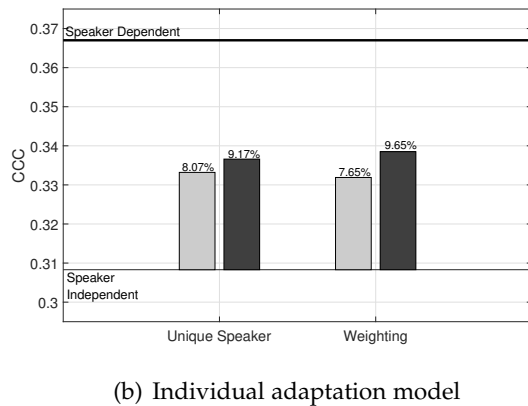
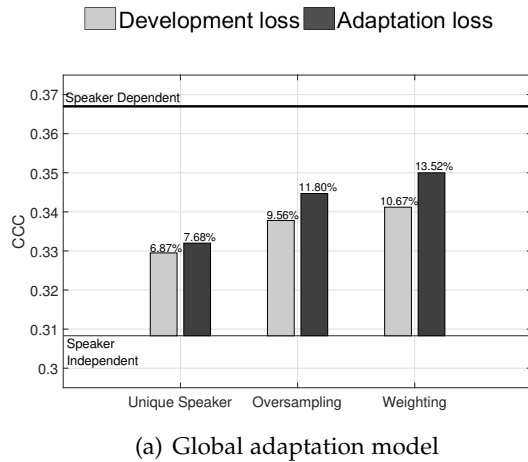


Fig. 10. Improvement in performance achieved by monitoring the loss function on the adaptation set while fine-tuning the models. The percentages over the bars indicate the relative improvements over the speaker-independent baseline. The figure shows the performance observed in the *test-B* set.

of the models. In our formulation, however, we aim to personalize a model towards a known set of speakers in the test set. Given the assumption that the selected speakers in the train set are similar to the test set, we can optimize the system by minimizing the loss function on the adaptation set when fine-tuning the model (i.e., samples from the selected speakers used for adaptation). Since we want to *personalize* the system, it is theoretically correct to maximize the performance of the model on data that is found to be closer to the target speakers. This section evaluates this idea. We record the best performance of the model by using an early stopping criterion on the adaptation loss. The only special case is the weighting approach which trains the models from scratch. Since monitoring the loss exclusively on the adaptation set will ignore other samples in the train set, we decide to monitor the loss on the full train set. The differences in the weights increase the emphasis of samples from the adaptation set, achieving essentially the same goal.

We use the adaptation set using the 200s condition. The weighting approach is implemented with the 1:3 ratio, which gave the best CCC in previous experiments (Figs. 5 and 7).

Figure 10 shows the performance improvements achieved by minimizing the loss function on the adaptation set. The darker bars indicate the results obtained while monitoring the training loss, and the lighter bars indicate the results obtained while monitoring the development loss. We include the relative improvements over the speaker-independent baseline with numbers on top of the bars. The relative improvements are constantly higher when maximizing the performance on the train set, tailoring even more the models to the test speakers. We conducted a one-tailed student t-test over ten trials to assert if the differences in performance between minimizing the loss function on the development and train sets are statistically significant. We assert significance when p -values < 0.05 . The statistical test indicates that the differences are statistically significant for all the adaptation approaches implemented with either the global or individual adaptation models. This approach leads to performances that are closer to the speaker-dependent baseline.

6.6 Performance on Other Emotional Attributes

The premise of this study is that valence is externalized in speech with more speaker-dependent traits than arousal and dominance. Therefore, personalization approaches to bring the models closer to the test speakers should have a higher impact on valence. This section implements the proposed adaptation schemes on arousal and dominance, comparing the relative improvements over the speaker-independent baseline with the results for valence.

Table 2 reports the performance and relative improvements over the speaker-independent baseline for valence, arousal, and dominance when using the proposed methods. The relative improvements for arousal and dominance are less than 1.9%, mirroring the results observed in Table 1 that shows relative improvements less than 3% for arousal and dominance when labeled data from the target speakers is available (i.e., speaker-dependent condition). Therefore, it is not surprising that the method does not lead to big improvements for arousal and dominance where there is little room for improvements. We argue that the approach is successful even for arousal or dominance, since the relative improvements for these emotional attributes are similar to the values reported in Table 1 under speaker dependent conditions. In contrast,

TABLE 2

Performance achieved using different adaptation approaches on the *test-B* set. The table reports the performance gain over the speaker-independent baseline reported in Table 1

	Adaptation Scheme	Minimizing loss in development set		Minimizing loss in adaptation set	
		CCC	Gain (%)	CCC	Gain (%)
Valence	Unique speaker (GA)	0.3295	6.87	0.3320	7.68
	Oversampling (GA)	0.3378	9.56	0.3447	11.80
	Weighting (GA)	0.3412	10.67	0.3500	13.52
	Unique speaker (IA)	0.3332	8.07	0.3366	9.17
	Weighting (IA)	0.3319	7.65	0.3385	9.79
Arousal	Unique speaker (GA)	0.7196	0.44	0.7221	0.79
	Oversampling (GA)	0.7209	0.62	0.7258	1.31
	Weighting (GA)	0.7222	0.80	0.7296	1.84
	Unique speaker (IA)	0.7185	0.29	0.7267	1.43
	Weighting (IA)	0.7202	0.53	0.7271	1.49
Dominance	Unique speaker (GA)	0.6410	0.56	0.6415	0.64
	Oversampling (GA)	0.6428	0.84	0.6430	0.87
	Weighting (GA)	0.6433	0.92	0.6451	1.20
	Unique speaker (IA)	0.6399	0.39	0.6419	0.70
	Weighting (IA)	0.6417	0.67	0.6422	0.75

TABLE 3

Comparison of results in terms CCC obtained with the proposed personalization approach and other methods. All the experiments are done with MSP-Podcast corpus and evaluated on Test-B set. STL: Single Task Learning (Speaker-Independent baseline), MTL: Multi-Task Learning

Attributes	Proposed Approach	STL	MTL	Ladder Networks
Valence	0.3500	0.3083	0.3302	0.3158
Arousal	0.7296	0.7164	0.7214	0.7421
Dominance	0.6451	0.6374	0.6287	0.6498

the relative improvements for valence are as high as 13.52%. These results validate our hypothesis that exploiting speaker-dependent characteristics between train and test speakers helps to personalize a SER system in the prediction of valence.

6.7 Comparison with Other Baselines

We compare the results from our proposed personalization approach with other state-of-the-art approaches. We consider *multi-task learning* (MTL) [13], and ladder network [14] for SER. These are some of the most successful approaches used in SER. The MTL approach jointly predicts arousal, valence and dominance, where the loss function is a weighted sum of the individual attribute losses. We used $\mathcal{L} = (1 - CCC)$ as the loss for each attribute (\mathcal{L}_{aro} , \mathcal{L}_{val} , \mathcal{L}_{dom}). Equation 5 shows the overall loss function, where $(\alpha, \beta) \in [0, 1]$ and $\alpha + \beta \leq 1$. The hyperparameters α and β are tuned on the development set.

$$\mathcal{L}_{MTL} = \alpha \mathcal{L}_{aro} + \beta \mathcal{L}_{val} + (1 - \alpha - \beta) \mathcal{L}_{dom} \quad (5)$$

The ladder network approach follows the implementation presented by Parthasarathy and Busso

TABLE 4

Comparison of CCC values between speaker-independent and speaker-dependent conditions for IEMOCAP (IEM) and MSP-IMPROV (IMP) databases. The DNN is trained with three layers and 256 nodes per layer. The column 'Gain' shows the relative improvement by training with partial data from the target speakers (*test-A* set).

Attributes	Database	Speaker Independent	Speaker Dependent	Gain (%)
		<i>Test-B</i> set	<i>Test-B</i> set	
Valence	IEM	0.4428	0.5072	14.54
	IMP	0.3420	0.4164	21.75
Arousal	IEM	0.6953	0.7255	4.34
	IMP	0.5958	0.6218	4.36
Dominance	IEM	0.5444	0.5678	4.29
	IMP	0.5625	0.4911	6.18

TABLE 5

IEMOCAP and MSP-IMPROV: Performance achieved using different adaptation approaches on the *test-B* set. The table reports the performance gain over the speaker-independent baselines reported in Table 4. All the experimental evaluations are done by minimizing the loss in the development set.

	Adaptation Scheme	USC-IEMOCAP		MSP-IMPROV	
		CCC	Gain (%)	CCC	Gain (%)
Valence	Unique speaker (GA)	0.4761	7.52	0.3866	13.04
	Oversampling (GA)	0.4790	8.17	0.4014	17.36
	Weighting (GA)	0.4889	10.41	0.3915	15.52
	Unique speaker (IA)	0.4725	6.70	0.3855	12.71
	Weighting (IA)	0.4759	7.47	0.3873	13.24
Arousal	Unique speaker (GA)	0.6988	0.50	0.6057	1.66
	Oversampling (GA)	0.7001	0.69	0.6086	2.14
	Weighting (GA)	0.7002	0.70	0.6191	3.91
	Unique speaker (IA)	0.6966	0.18	0.6087	2.16
	Weighting (IA)	0.7000	0.68	0.6095	2.29
Dominance	Unique speaker (GA)	0.5491	0.86	0.4718	2.01
	Oversampling (GA)	0.5500	1.02	0.4731	2.29
	Weighting (GA)	0.5495	0.93	0.4820	4.21
	Unique speaker (IA)	0.5451	0.12	0.4710	1.83
	Weighting (IA)	0.5496	0.95	0.4713	1.90

[14]. This method uses the reconstruction of feature representations at various layers in a DNN as auxiliary tasks. In addition, we consider the speaker-independent baseline discussed in previous sections, referred here to as *single task learning* (STL).

Table 3 describes the results, which clearly show significant improvements over alternative methods in CCC for valence using our proposed approach, reinforcing our claim about the benefits of personalization in the estimation of valence. Other approaches are effective for arousal and dominance.

6.8 Performance on Other Corpora

We validate the effectiveness of the proposed personalization approaches with other emotional databases. We use a K -fold cross-validation strategy to train DNN models using the USC-IEMOCAP and MSP-IMPROV databases. In each fold, we consider 2 speakers as the test speakers and the rest of

the speakers as the train speakers. With this cross-validation approach, all the speakers are at some point considered as the test speakers. The final results are averaged across the K folds. Similar to Figure 1, we split the test set of both IEMOCAP and MSP-IMPROV databases into two partitions for this study: *test-A* and *test-B* sets. The *test-A* set includes 200s of recording for each of the speakers in the test set. This set is reserved for finding the closest training speakers to the target speakers. The *test-B* set includes the rest of the recordings in the test set. We implement the global and individual adaptation models with all the different adaptation approaches by minimizing the loss in the development set.

Table 4 shows the CCC results for speaker-independent and speaker-dependent conditions using the USC-IEMOCAP and the MSP-IMPROV corpora. The results validate the findings in Table 1, showing that the speaker-dependent condition leads to higher relative gains for valence than for arousal or dominance.

Table 5 shows the results obtained with different adaptation schemes on the USC-IEMOCAP and MSP-IMPROV databases. The results are consistent with the findings observed with the MSP-Podcast corpus. We observe that the relative improvements over the speaker-independent baseline are much higher for valence than the relative improvements for arousal and dominance. With the USC-IEMOCAP corpus, we achieve relative gains in performance up to 10.41% for valence whereas less than 1.03% for arousal and dominance. Similarly, with the MSP-IMPROV corpus, we achieve relative gains in performance up to 17.36% for valence whereas less than 4.22% for arousal and dominance. These results show the effectiveness of our proposed approach applied to other emotional corpora, reinforcing our finding about the speaker-dependent nature of valence emotional cues.

7 CONCLUSIONS

This paper demonstrated that a valence prediction system can be personalized to target speakers by exploiting speaker-dependent traits. The study proposed to create an adaptation set by identifying speakers in the train set that are closer to the speakers in the test set. Since we evaluate the similarity between speakers by comparing the acoustic feature spaces associated with each speaker without using emotional labels, the adaptation approaches are fully unsupervised. We proposed three methods to create this adaptation set: unique speaker,

oversampling and weighting approaches. The adaptation sets from the selected *closer* speakers are used to personalize the DNNs to the speakers in the test set. The experimental results showed that the global adaptation models achieved better performance than the individual adaptation models. Further improvements are observed when the loss function is minimized by monitoring the loss in the adaptation set. The proposed adaptation schemes lead to relative improvements up to 13.52% over a speaker-independent baseline. We observed significant improvements in performance, even when only a few seconds of adaptation data (belonging to the train set) for each of the speakers in the test set was used for adaptation. The maximum performance gains were observed for 200s of adaptation data for each of the speakers in the test set. We also demonstrated the effectiveness of the proposed personalization approaches with the USC-IEMOCAP and MSP-IMPROV databases, showing consistent findings.

There are many interesting and important applications where our personalization models can be very helpful. In healthcare applications where medical data cannot be frequently obtained, a personalized system can keep track of a patient's expressive behaviors and his/her medical record for better and efficient treatment. Another example is on personal assistant systems on mobile devices, which are often used by a limited number of users. Over time, the system can collect enough data from the target users to improve their emotion recognition systems. For cloud-based applications, the training data needs to be stored in the cloud, instead of the edge device. With a simple modification on the approach to obtain the PCA projections (i.e., finding a common PCA space across testing speakers), the training data can be pre-estimated and stored. Therefore, this approach does not require storage or computational resources on the edge devices and can be efficiently implemented during inferences.

This study demonstrated the importance of exploiting the speaker-dependent traits observed in the expression of valence from speech, which led to clear improvements. It also showed that we can personalize SER models by just finding speakers in the train set that are similar to the target speakers. The proposed formulation is flexible, requiring only to know the block of data associated with each speaker in the test set. This assumption is reasonable since it is straightforward to group data per speaker in many practical applications. As a future work, we will evaluate more sophisticated methods to assess

the similarity of speakers by considering more than acoustic similarities. Also, we will explore the use of the proposed adaptation schemes in other deep learning frameworks such as autoencoders, *generative adversarial networks* (GANs) or *long short-term memory* (LSTM).

ACKNOWLEDGMENT

This study was funded by National Science Foundation (NSF) career award IIS-1453781

REFERENCES

- [1] Z. Aldeneh, M. Jaiswal, M. Picheny, M. McInnis, and E. Mower Provost, "Identifying mood episodes using dialogue features from clinical interviews," in *Interspeech 2019*, Graz, Austria, September 2019, pp. 1926–1930.
- [2] J. Gideon, H. Schatten, M. McInnis, and E. Mower Provost, "Emotion recognition from natural phone conversations in individuals with and without recent suicidal ideation," in *Interspeech 2019*, Graz, Austria, September 2019, pp. 3282–3286.
- [3] S. Narayanan and P. G. Georgiou, "Behavioral signal processing: Deriving human behavioral informatics from speech and language," *Proceedings of the IEEE*, vol. 101, no. 5, pp. 1203–1233, May 2013.
- [4] P. Georgiou, M. Black, and S. Narayanan, "Behavioral signal processing for understanding (distressed) dyadic interactions: some recent developments," in *joint ACM workshop on Human gesture and behavior understanding (J-HGBU 2011)*, Scottsdale, Arizona, USA, December 2011, pp. 7–12.
- [5] J. A. Russell, "Core affect and the psychological construction of emotion." *Psychological review*, vol. 110, no. 1, p. 145, 2003.
- [6] N. A. Groenewold, E. M. Opmeer, P. de Jonge, A. Aleman, and S. G. Costafreda, "Emotional valence modulates brain functional abnormalities in depression: evidence from a meta-analysis of fmri studies," *Neuroscience & Biobehavioral Reviews*, vol. 37, no. 2, pp. 152–163, 2013.
- [7] M. Hagenhoff, N. Franzen, L. Gerstner, G. Koppe, G. Sammer, P. Netter, B. Gallhofer, and S. Lis, "Reduced sensitivity to emotional facial expressions in borderline personality disorder: effects of emotional valence and intensity," *Journal of personality disorders*, vol. 27, no. 1, pp. 19–35, 2013.
- [8] M. Nicolaou, H. Gunes, and M. Pantic, "Continuous prediction of spontaneous affect from multiple cues and modalities in valence-arousal space," *IEEE Transactions on Affective Computing*, vol. 2, no. 2, pp. 92–105, April-June 2011.
- [9] Z. Aldeneh, S. Khorram, D. Dimitriadis, and E. Mower Provost, "Pooling acoustic and lexical features for the prediction of valence," in *ACM International Conference on Multimodal Interaction (ICMI 2017)*, Glasgow, UK, November 2017, pp. 68–72.
- [10] B. Zhang, S. Khorram, and E. Mower Provost, "Exploiting acoustic and lexical properties of phonemes to recognize valence from speech," in *IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP 2019)*, Brighton, United Kingdom, May 2019, pp. 5871–5875.
- [11] E. Tournier, "Valence of emotional memories: A study of lexical and acoustic features in older adult affective speech," Master's thesis, Radboud University Nijmegen, Nijmegen, Netherlands, June 2019.
- [12] S. Mariooryad and C. Busso, "Exploring cross-modality affective reactions for audiovisual emotion recognition," *IEEE Transactions on Affective Computing*, vol. 4, no. 2, pp. 183–196, April-June 2013.
- [13] S. Parthasarathy and C. Busso, "Jointly predicting arousal, valence and dominance with multi-task learning," in *Interspeech 2017*, Stockholm, Sweden, August 2017, pp. 1103–1107.
- [14] —, "Ladder networks for emotion recognition: Using unsupervised auxiliary tasks to improve predictions of emotional attributes," in *Interspeech 2018*, Hyderabad, India, September 2018, pp. 3698–3702.
- [15] K. Sridhar, S. Parthasarathy, and C. Busso, "Role of regularization in the prediction of valence from speech," in *Interspeech 2018*, Hyderabad, India, September 2018.
- [16] G. Trigeorgis, F. Ringeval, R. Brueckner, E. Marchi, M. Nicolaou, B. Schuller, and S. Zafeiriou, "Adieu 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)*, Shanghai, China, March 2016, pp. 5200–5204.
- [17] C. Busso and T. Rahman, "Unveiling the acoustic properties that describe the valence dimension," in *Interspeech 2012*, Portland, OR, USA, September 2012, pp. 1179–1182.
- [18] S. Koolagudi and K. Rao, "Exploring speech features for classifying emotions along valence dimension," in *International Conference on Pattern Recognition and Machine Intelligence (PReMI 2019)*, ser. Lecture Notes in Computer Science, S. Chaudhury, S. Mitra, C. Murthy, P. Sastry, and S. Pal, Eds. New Delhi, India: Springer, Berlin, Heidelberg, December 2009, vol. 5909, pp. 537–542.
- [19] N. Cook, T. Fujisawa, and K. Takami, "Evaluation of the affective valence of speech using pitch substructure," *IEEE Transactions on Audio, Speech, and Language Processing*, vol. 14, no. 1, pp. 142–151, January 2005.
- [20] N. Cook and T. Fujisawa, "The use of multi-pitch patterns for evaluating the positive and negative valence of emotional speech," in *Speech Prosody (SP 2006)*, Dresden, Germany, May 2006, pp. 5–8.
- [21] G. Deshpande, V. S. Viraraghavan, M. Duggirala, and S. Patel, "Detecting emotional valence using time-domain analysis of speech signals," in *IEEE International Engineering in Medicine and Biology Conference (EMBC 2019)*, Berlin, Germany, July 2019, pp. 3605–3608.
- [22] G. Deshpande, V. Viraraghavan, and R. Gavas, "A successive difference feature for detecting emotional valence from speech," in *Workshop on Speech, Music and Mind (SMM 2019)*, Vienna, Austria, September 2019, pp. 36–40.
- [23] A. Tursunov, S. Kwon, and H. Pang, "Discriminating emotions in the valence dimension from speech using timbre features," *Applied Sciences*, vol. 9, no. 12, p. 18, June 2019.
- [24] M. Tahon, G. Degottex, and L. Devillers, "Usual voice quality features and glottal features for emotional valence detection," in *Speech Prosody (SP 2012)*, Shanghai, China, May 2012, pp. 693–696.
- [25] C.-C. Lee, C. Busso, S. Lee, and S. Narayanan, "Modeling mutual influence of interlocutor emotion states in dyadic spoken interactions," in *Interspeech 2009*, Brighton, UK, September 2009, pp. 1983–1986.
- [26] E. Lakomkin, M. Zamani, C. Weber, S. Magg, and S. Wermter, "Incorporating end-to-end speech recognition

- models for sentiment analysis," in *International Conference on Robotics and Automation (ICRA 2019)*, Montreal, QC, Canada, May 2019, pp. 7976–7982.
- [27] M. Abdelwahab and C. Busso, "Supervised domain adaptation for emotion recognition from speech," in *International Conference on Acoustics, Speech, and Signal Processing (ICASSP 2015)*, Brisbane, Australia, April 2015, pp. 5058–5062.
- [28] —, "Incremental adaptation using active learning for acoustic emotion recognition," in *IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP 2017)*, New Orleans, LA, USA, March 2017, pp. 5160–5164.
- [29] Y. Bengio, "Deep learning of representations for unsupervised and transfer learning," in *ICML Workshop on Unsupervised and Transfer Learning*, Bellevue, WA, USA, July 2012, pp. 17–36.
- [30] J. Deng, Z. Zhang, E. Marchi, and B. Schuller, "Sparse autoencoder-based feature transfer learning for speech emotion recognition," in *2013 Humaine Association Conference on Affective Computing and Intelligent Interaction*. IEEE, 2013, pp. 511–516.
- [31] 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 *Proc. International Conference on Acoustics, Speech, and Signal Processing (ICASSP)*, 2014, pp. 4851–4855.
- [32] J. Gideon, S. Khorram, Z. Aldeneh, D. Dimitriadis, and E. M. Provost, "Progressive neural networks for transfer learning in emotion recognition," *arXiv preprint arXiv:1706.03256*, 2017.
- [33] A. Hassan, R. Damper, and M. Niranjana, "On acoustic emotion recognition: compensating for covariate shift," *IEEE Transactions on Audio, Speech, and Language Processing*, vol. 21, no. 7, pp. 1458–1468, 2013.
- [34] P. Song, S. Ou, W. Zheng, Y. Jin, and L. Zhao, "Speech emotion recognition using transfer non-negative matrix factorization," in *Acoustics, Speech and Signal Processing (ICASSP), 2016 IEEE International Conference on*. IEEE, 2016, pp. 5180–5184.
- [35] 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.
- [36] Z. Huang, W. Xue, Q. Mao, and Y. Zhan, "Unsupervised domain adaptation for speech emotion recognition using pcanet," *Multimedia Tools and Applications*, vol. 76, no. 5, pp. 6785–6799, 2017.
- [37] T. Rahman and C. Busso, "A personalized emotion recognition system using an unsupervised feature adaptation scheme," in *International Conference on Acoustics, Speech, and Signal Processing (ICASSP 2012)*, Kyoto, Japan, March 2012, pp. 5117–5120.
- [38] C. Busso, S. Mariooryad, A. Metallinou, and S. Narayanan, "Iterative feature normalization scheme for automatic emotion detection from speech," *IEEE Transactions on Affective Computing*, vol. 4, no. 4, pp. 386–397, October–December 2013.
- [39] R. Lotfian and C. Busso, "Building naturalistic emotionally balanced speech corpus by retrieving emotional speech from existing podcast recordings," *IEEE Transactions on Affective Computing*, vol. 10, no. 4, pp. 471–483, October–December 2019.
- [40] S. Mariooryad, R. Lotfian, and C. Busso, "Building a naturalistic emotional speech corpus by retrieving expressive behaviors from existing speech corpora," in *Interspeech 2014*, Singapore, September 2014, pp. 238–242.
- [41] A. Burmania, S. Parthasarathy, and C. Busso, "Increasing the reliability of crowdsourcing evaluations using online quality assessment," *IEEE Transactions on Affective Computing*, vol. 7, no. 4, pp. 374–388, October–December 2016.
- [42] C. Busso, M. Bulut, C. Lee, A. Kazemzadeh, E. Mower, S. Kim, J. Chang, S. Lee, and S. Narayanan, "IEMOCAP: Interactive emotional dyadic motion capture database," *Journal of Language Resources and Evaluation*, vol. 42, no. 4, pp. 335–359, December 2008.
- [43] C. Busso, S. Parthasarathy, A. Burmania, M. AbdelWahab, N. Sadoughi, and E. Mower Provost, "MSP-IMPROV: An acted corpus of dyadic interactions to study emotion perception," *IEEE Transactions on Affective Computing*, vol. 8, no. 1, pp. 67–80, January–March 2017.
- [44] B. Schuller, S. Steidl, A. Batliner, A. Vinciarelli, K. Scherer, F. Ringeval, M. Chetouani, F. Wenginger, F. Eyben, E. Marchi, M. Mortillaro, H. Salamin, A. Polychroniou, F. Valente, and S. Kim, "The INTERSPEECH 2013 computational paralinguistics challenge: Social signals, conflict, emotion, autism," in *Interspeech 2013*, Lyon, France, August 2013, pp. 148–152.
- [45] N. Dehak, P. A. Torres-Carrasquillo, D. Reynolds, and R. Dehak, "Language recognition via i-vectors and dimensionality reduction," in *Twelfth annual conference of the international speech communication association*. Citeseer, 2011.
- [46] D. Snyder, D. Garcia-Romero, G. Sell, D. Povey, and S. Khudanpur, "X-vectors: Robust dnn embeddings for speaker recognition," in *2018 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 2018, pp. 5329–5333.

Kusha Sridhar (S'18) received his BS in Electronics and Communications Engineering from PES University, Bangalore, Karnataka, India, in 2015 and MS degree in Electrical Engineering from the University of Southern California (USC), Los Angeles, in 2017. He is currently pursuing his Ph.D. degree in Electrical Engineering at the University of Texas at Dallas. His research interests include areas related to affective computing, focusing on emotion recognition from speech, machine learning and speech signal processing.



Carlos Busso (S'02-M'09-SM'13) is a professor in the Electrical Engineering Department of the University of Texas at Dallas, Dallas, Texas, 75080, USA. He received his Ph.D. degree in electrical engineering from the University of Southern California, Los Angeles, California, USA. He is a recipient of an NSF CAREER Award. In 2014, he received the ICMI Ten-Year Technical Impact Award. He was a core-cipient of both the Hewlett-Packard Best Paper Award at the 2011 IEEE International Conference on Multimedia and Expo and the Best Paper Award at the Association for the Advancement of Affective Computing International Conference on Affective Computing and Intelligent Interaction 2017. His research interest include human-centered multimodal machine intelligence and applications; and the broad areas of affective computing, nonverbal behaviors for conversational agents, and machine learning methods for multimodal processing.

