# Multi-Task Semi-Supervised Adversarial Autoencoding for Speech Emotion Recognition

Siddique Latif, Rajib Rana, Sara Khalifa, Raja Jurdak, Julien Epps, and Björn W. Schuller, *Fellow, IEEE*

**Abstract**—Inspite the emerging importance of Speech Emotion Recognition (SER), the state-of-the-art accuracy is quite low and needs improvement to make commercial applications of SER viable. A key underlying reason for the low accuracy is the scarcity of emotion datasets, which is a challenge for developing any robust machine learning model in general. In this paper, we propose a solution to this problem: a multi-task learning framework that uses auxiliary tasks for which data is abundantly available. We show that utilisation of this additional data can improve the primary task of SER for which only limited labelled data is available. In particular, we use gender identifications and speaker recognition as auxiliary tasks, which allow the use of very large datasets, e. g., speaker classification datasets. To maximise the benefit of multi-task learning, we further use an adversarial autoencoder (AAE) within our framework, which has a strong capability to learn powerful and discriminative features. Furthermore, the unsupervised AAE in combination with the supervised classification networks enables semi-supervised learning which incorporates a discriminative component in the AAE unsupervised training pipeline. This semi-supervised learning essentially helps to improve generalisation of our framework and thus leads to improvements in SER performance. The proposed model is rigorously evaluated for categorical and dimensional emotion, and cross-corpus scenarios. Experimental results demonstrate that the proposed model achieves state-of-the-art performance on two publicly available datasets.

**Index Terms**—Speech emotion recognition, multi task learning, representation learning

✦

## 1 INTRODUCTION

SPEECH Emotion Recognition (SER) is an emerging area of research. Since speech is a major form of affect display [1], the success of SER will redefine human-computer interactions, enabling, for example, effective service delivery in many sectors. Call centres now track customers' emotions for better service delivery [2]. Speech based diagnostic systems are being developed for diagnosis of depression [3], distress [4], and monitoring of mood states for bipolar patients [5]. Many other applications including media retrieval systems [6], smart cars [7], and forensic sciences [8] also aim to improve their performances by utilising SER techniques.

Human emotions in speech are complex to model due to dependency of speech on many factors including speaker [9], gender [10], age [11], culture [12], dialect [13], and among others. Researchers have explored many methods including classical models, such as hidden Markov models, support vector classification, and deep neural networks (DNNs) for speech emotion recognition, wherein DNN models have usually demonstrated better performance compared to the classical models [14], [15]. Currently,

the popularity of DNN models for speech emotion recognition is seeing a steep rise.

DNN models that have been successful for speech emotion recognition include deep belief networks (DBN) [16], [17], convolutional neural networks (CNN) [18], [19] and long short term memory (LSTM) networks [20], [21], [22]. The majority of the above research presents techniques to predict speech emotion using single task (emotion recognition) training. These techniques, however, ignore a potentially rich source of information available in speech (e. g., information about the speaker, gender, etc.) that can be utilised for achieving generalisation and improvement in the performance [23]. To achieve generalisation, most existing studies tend to validate/tune models using diverse datasets [12], [17]. However, standard benchmark datasets are very scarce, and most problematically, they are of smaller sizes, which creates massive roadblocks in achieving generalisation in SER systems [23].

An alternative and effective approach to increase the generalisation of SER models is multi-task learning (MTL) [24], which simultaneously solves relevant auxiliary tasks along with the primary task. In MTL, models are better regularised to uncover the common high-level discriminative representations. MTL has been widely applied to various speech and natural language processing related problems [25], [26]. In SER, MTL has shown good performance for fully supervised deep learning (DL) models [23], [27], [28]. Most of these approaches jointly learn different emotional attributes to improve both performance and generalisation [28], [29], [30]. Often, researchers use categorical emotion as a primary task and dimensional emotion as an auxiliary task. Discrete/categorical theories of emotions encompass a small set of distinct emotions. The foundation

- S. Latif is affiliated with USQ, Australia and Distributed Sensing Systems Group, Data61, CSIRO Australia.
- R. Rana is with University of Southern Queensland (USQ), Australia.
- S. Khalifa is affiliated with Distributed Sensing Systems Group, Data61, CSIRO Australia.
- R. Jurdak is affiliated with Queensland University of Technology (QUT), Australia.
- J. Epps is affiliated with University of New South Wales (UNSW), Australia.
- B. Schuller is affiliated with GLAM – the Group on Language, Audio, and Music, Imperial College London, UK, and the ZD.B Chair of Embedded Intelligence for Health Care and Wellbeing, University of Augsburg, Germany.
Corresponding E-mail: siddique.latif@usq.edu.au

of these theories is that different emotions are associated with distinctive patterns of triggers, behavioural expression, and unique subjective experiences [31]. Only core emotions including joy, sadness, fear, anger, and disgust [32] are included in these theories. On the contrary, the foundation of the dimensional models of emotions is that a common and interconnected neurophysiological system generates all affective states [33]. Generally, these models define human emotions using a two-dimensional space having valence in one dimension and activation or intensity in the other dimension. To use these dimensional emotions as secondary tasks, annotation is important. However, the meta labels for these emotional attributes are scarcely available. Recently, it has been shown in computer vision that the performance of primary tasks with constrained data can be significantly enhanced by using larger data for the auxiliary tasks [34], [35]. Inspired by this idea, in this study, we aim to build models that can effectively utilise auxiliary tasks with a large quantity of available data in order to improve the performance of the primary task. We use emotion recognition as our primary task and select gender and speaker recognition as auxiliary tasks to include larger datasets.

Within our MTL framework, we further utilise generative adversarial models due to their exceptional ability to learn powerful and discriminative features [36]. In particular, we use adversarial autoencoder (AAE) [37], which fundamentally aims to learn representation of data in an unsupervised way. However, by combining AAE with the supervised classification networks we enable semi-supervised learning for AAE. This essentially incorporates the discriminative component (from the supervised classification) in the training pipeline of unsupervised learning to influence the latent representation of AAE and by makes it suitable for semi-supervised emotion classification.

To show the advantage of our proposed MTL framework, we evaluate it comprehensively on two large and widely used emotional databases: The interactive emotional dyadic motion capture (IEMOCAP) [38] and MSP-IMPROV [39]. We compare the performance of our proposed framework with that of recent studies, and also with popular models like CNNs, and an autoencoder based semi-supervised model. The comparative results show that, for categorical, dimensional, and cross-corpus emotion classification, we achieve the improved results, which attests to the strong generalisation power of the proposed framework.

## 2 RELATED WORK

Our framework utilises multi-task learning for SER. It also uses semi supervised learning while employing adversarial encoding, where the classification is done through CNNs. We therefore cover these four aspects in our literature review.

### 2.1 Landscape of Multi-task Learning for SER

Multi-task learning (MTL) has been successful for simultaneously modelling multiple related tasks utilising shared representation [40], [41]. It aims to improve generalisation by learning the similarities as well as the differences among the given tasks from the training data [24]. The conventional

methodology to optimise a machine learning model for one task at a time ignores potentially rich information in the training signal [23]. Such information can be effectively utilised for auxiliary tasks to improve generalisation and performance of a system. Several MTL approaches [42], [43], [44] have been widely used for solving problems in computer vision. The primary reason to use MTL in vision is that images can provide information related to different tasks, and simultaneously learning these correlated tasks can boost the performance of each individual task [45], [46]. For example, face detection, gender recognition, and pose estimation can be simultaneously performed using a deep CNN [47].

Similarly to images, speech is another such modality that can provide information for various tasks including speaker, gender, and emotion identification. Researchers have started to investigate the effectiveness of MTL for improving the performance of speech emotion recognition [27], [48], [49]. Eyben et al. [50] were first to use MTL in SER and they showed that training the model with multiple targets helps to improve the performance compared to single target training. Prthasarathy and Busso [28] proposed a DNN based model to jointly learn the arousal, dominance, and valence value of a given utterance. The authors demonstrated that joint learning of these emotional attributes significantly enhances the performance of a model compared to single task learning (STL). Similarly, Ma et al. [49] used a multi-task attention-based DNN for SER and found that, by sharing the information among tasks, a high performance can be achieved. Xia et al. [30] proposed a DBN based model for MTL and utilised activation and valence information for speech emotion recognition. They illustrated that the utilisation of additional information in the MTL setup can improve the performance of their model by considering the categorical emotion label as the primary task, and activation and valence information as secondary tasks. Similarly, Lotfian et al. [29] used a DNN for jointly learning primary and secondary emotions. They showed that the classification performance of the primary task (categorical emotions) is significantly improved by considering secondary emotions (other emotional classes perceived by the evaluators) in the model. In another study, Chang et al. [51] used a generative adversarial network (GAN) for MTL with valence classification as primary and activation classification as a secondary task. In addition, the authors also introduced unlabelled data from the AMI corpus [52] (a multi-modal data set consisting of 100 hours of meeting recordings) to train generator and discriminator components of a GAN and showed that the performance of the classifier can be improved by using task-unrelated speech data in an unsupervised way.

Another stream of research in SER—instead of using different emotional attributes as auxiliary tasks—has utilised other available attributes, such as speaker identity and gender to improve the performance of SER [53]. For instance, Kim et al. [27] used gender and naturalness (natural or acted corpus) recognition as auxiliary tasks to improve the performance of emotion recognition using different emotional databases. Zhang et al. [54] used an MTL approach to investigate the influence of the domain (whether the expression is spoken or sung), corpus, and gender on cross-corpus emotion recognition systems. The authors used six different

emotional databases and showed that the performance of a cross-corpus SER system increases with the rising number of emotional corpora used for training. Based on these results, they also showed that effective modelling of cross-corpus emotion recognition requires the understanding of emotional changes as a function of non-emotional factors.

Both streams of research mentioned above conform to the fact that MTL approaches can improve the SER performance compared with STL. While the first stream shows that choosing emotional attributes as auxiliary tasks leads to improved performance of SER for the primary task, the second stream shows that it is also possible to choose non-emotional attributes of speech as a secondary task, and performance of SER as a primary task can be improved. Our approach is motivated by the second stream as it provides the opportunity to utilise abundantly available non-emotional datasets. Precisely, we consider using abundantly available non-emotional speech corpora to indirectly improve the performance of SER by directly improving the performance of the auxiliary tasks, which has not been widely studied in the existing literature. In [51], the authors used additional data, however, unlike them, we use additional data for auxiliary tasks. Also, unlike them, we backpropagate AE reconstruction loss in addition to backpropagating classification, generator, and discriminator losses. Note that, as our training uses both labelled and unlabelled emotion data, we therefore, introduce semi-supervised learning in MTL. In the next section, we cover studies using semi-supervised learning for SER.

### 2.2 Landscape of Semi-Supervised Learning for SER

A number of studies have considered semi-supervised learning for SER. Huang et al. [60] introduced semi-supervised CNN for learning affect-salient features and reported superior performance on four public emotional speech databases: the Surrey Audio-Visual Expressed Emotion (SAVEE) database [61], the Berlin Emotional Database (Emo-DB) [62], the Danish Emotional Speech database (DES) [63], and the Mandarin Emotional Speech database (MES) [64]. The authors used CNN in an unsupervised way to learn general features and then fine-tuned the model for emotion recognition. Zhang et al. [55] proposed a collaborative semi-supervised learning technique that can correct mislabelled samples by re-evaluating the automatically labelled samples in learning iterations of the model. They also used different models including SVMs and RNNs and multiple modalities (audio and video) to improve the performance by simultaneously minimising the joint entropy. Recently, researchers further studied ladder network-based semi-supervised methods for SER [57], [58], [59], [65] and have shown superior results over supervised methods. A ladder network is an unsupervised denoising autoencoder that is trained along with a supervised classification or regression task. Deng et al. [56] proposed a framework for SER by combining an autoencoder and a classifier. Their work is based on a discriminative Restricted Boltzmann Machine (RBM) [66], which considers unlabelled samples as an extra garbage class in the classification problem. Our study differs from previous studies by simultaneously training an adversarial autoencoder with multi-task classifiers and utilising

the additional unlabelled emotional data for auxiliary tasks to improve SER performance. Joint optimisation of the sum of multi-task supervised and unsupervised cost functions is an important contribution leading to more discriminative SER models. In the next section, we focus on the existing studies that utilise adversarial autoencoders (AAE) for SER and highlight the difference with our work.

### 2.3 Landscape of Adversarial Autoencoders (AAE) for SER

Autoencoders are unsupervised learning models that have been successfully utilised in the field of SER. They are very powerful in learning salient representations that lead to a notable improvement in SER performance [20], [67]. Adversarial autoencoders (AAEs) are probabilistic models [37] that turn an autoencoder into a generative model. This has increased the popularity of AAEs in learning more descriptive features compared to conventional AEs and even compared to variational autoencoders (VAEs) [68]. In [36], AAEs have been used in SER for encoding high dimensional feature representations into compressed space and for the generation of speech samples. The authors found that the latent code learnt by AAE preserves class discriminability that is very crucial for speech emotion classification. However, most SER studies utilised AE networks to perform feature learning and then classification was performed separately (e. g., [36], [69]). However, it has been shown in that AAEs can be exploited in semi-supervised way to improve the classification performance [37]. Therefore, we proposed a self-sufficient semi-supervised structure that can performs both feature representation learning and classification learning by jointly minimising reconstruction error and the sum of multi-task classification errors.

### 2.4 Landscape of CNNs for SER

Convolutional neural networks (CNNs) are one of the most popular deep learning models that have demonstrated great success in various research fields including object recognition [70], handwriting recognition [71], face recognition [72], natural language processing (NLP) [73], and speech recognition [74]. CNNs overcome the scalability problem of standard neural networks by allowing the multiple regions of the input to share the same weights [15]. Generally, CNNs consist of three building blocks: convolutional layers, pooling layers, and fully connected layers. Convolutional layers in CNNs perform a convolution operation to compute feature maps, which are then sub-sampled using pooling layers. Finally, fully connected layers are used to transform the features into a more discriminative space for target prediction. In SER, CNNs have been widely used to learn salient features [18], [75], [76], also directly for classification [77]. Studies [78], [79], [80], [81] also presented CNNs in combination with LSTM to improve SER performance. However, this study proposes a unique use of CNN upon using it in an MTL framework while utilising the unlabelled data for the auxiliary task to improve SER performance. For the convenience of the readers, in Table 1 we provide a difference of our work with that of the existing literature, which supports the claims we make in this paper.

Table 1: Summary of comparative analysis of our paper with that of the existing literature.

| Paper/Author (Year) | Auxiliary Tasks for MTL | | Adversarial Learning | Semi-Supervised Learning | Additional Data for Auxiliary Tasks |
| --- | --- | --- | --- | --- | --- |
| | Emotional Attributes | Non-Emotional Attributes | | | |
| Prthasarathy and Busso [28] (2017) | ✓ | ✗ | ✗ | ✗ | ✗ |
| Xia et al. [30] (2017) | ✓ | ✗ | ✗ | ✗ | ✗ |
| Chang et al. [51] (2017) | ✓ | ✗ | ✓ | ✓ | ✗ |
| Lotfian et al. [29] (2018) | ✓ | ✗ | ✗ | ✗ | ✗ |
| Tao et al. [53] (2018) | ✗ | ✓ | ✗ | ✗ | ✗ |
| Zhang et al. [55] (2018) | ✗ | ✗ | ✗ | ✓ | ✗ |
| Deng et al. [56] (2018) | ✗ | ✗ | ✗ | ✓ | ✗ |
| Huang et al. [57] (2018) | ✗ | ✗ | ✗ | ✓ | ✗ |
| Sahu et al. [36] (2018) | ✗ | ✗ | ✓ | ✗ | ✗ |
| Tao et al. [58] (2019) | ✗ | ✗ | ✗ | ✓ | ✗ |
| Prthasarathy and Busso [59] (2019) | ✓ | ✗ | ✗ | ✓ | ✗ |
| **Our Paper (2019)** | ✗ | ✓ | ✓ | ✓ | ✓ |

## 3 PROPOSED MODEL

We proposed a multi-task learning framework by incorporating semi-supervised adversarial autoencoding using adversarial autoencoders (AAE). An AAE combines a traditional autoencoder and an adversarial network to deliver a surprisingly flexible framework. In AAE, the adversarial part is attached to the latent code z, where the encoder of autoencoder network also acts as the generator of the adversarial network. It enforces the autoencoder to generate a latent representation $z$ by observing the statistical properties of a given prior distribution $p(d)$.
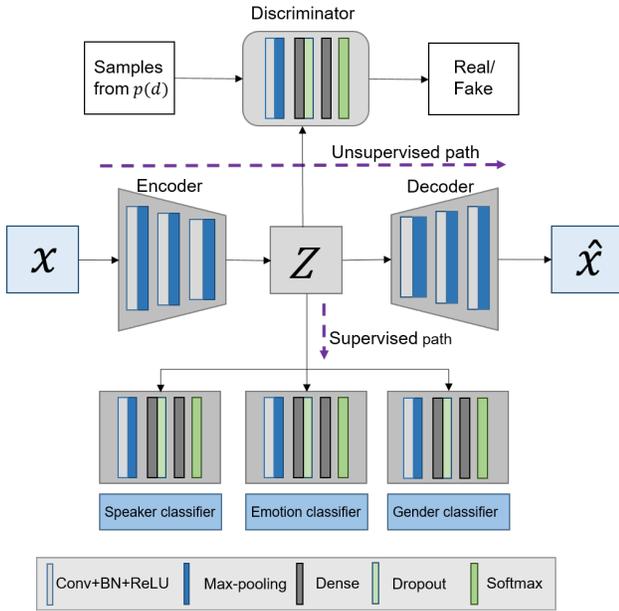


Figure 1: Illustration of our proposed multitask framework using a semi-supervised adversarial autoencoder (AAE).

In order to achieve MTL in AAE, we modify it to incorporate three supervised classification networks including emotion, speaker, and gender classification. Fig. 1 illustrates our proposed semi-supervised multitask learning model, where we highlight the supervised and unsupervised paths. In Equation (1) we present the multi-task autoencoding loss $\mathcal{L}_{\text{MTAE}}$ as a function of supervised and unsupervised losses.

$$\mathcal{L}_{\text{MTAE}} = \alpha * \mathcal{L}_{\text{AE}} + \mathcal{L}_c, \quad (1)$$

$$\mathcal{L}_c = \beta * \mathcal{L}_{\text{E}} + (1 - \beta) * (\mathcal{L}_{\text{G}} + \mathcal{L}_{\text{S}}). \quad (2)$$

Here, $\mathcal{L}_{\text{AE}}$ is the reconstruction loss of the autoencoder; $\mathcal{L}_{\text{E}}$, $\mathcal{L}_{\text{G}}$, and $\mathcal{L}_{\text{S}}$ are losses for the emotion, gender, and speaker classification tasks, respectively; $\alpha$ and $\beta$ are the trade-off parameters to control the weight of each loss term. In addition to the autoencoding network (encoder ($E_\theta$) and decoder ($D_\delta$)) and the and classifiers ($C_\phi$), there is an adversarial network that includes a generator ($E_\theta$) and discriminator ($D_\omega$).

For the input data $x$, overall model is trained in three phases: (1) the reconstruction phase; (2) the regularisation phase; and (3) the classification phase. In the reconstruction phase, the autoencoder updates the encoder ($E_\theta$) and the decoder ($D_\delta$) and minimises the reconstruction error by encoding $x$ into latent representation $z$. The objective function for the autoencoder is defined below:

$$\mathcal{L}_{\text{AE}}(x, D_\delta(E_\theta(x))) = \|x - \hat{x}\|_2^2 \quad (3)$$

In *the regularisation phase*, the adversarial network first updates its discriminator ($D_\omega$) to distinguish between the samples coming from the prior distribution $p(d) = \mathcal{N}(d; 0, I)$ [real] and that generated using the latent codes ($z$) [fake] computed by the autoencoder; and then updates its generator ($E_\theta$) or encoder . The objective here is to fool the discriminator ($D_\omega$) by learning to encode data that $D_\omega$ perceives as real. The update is done by keeping the weight and bias of the discriminator network fixed and by backpropagating the error to $E_\theta$ and updating its weight and bias values. The objective function for the discriminator ($D_\omega$) is defined as:

$$L_{\text{disc}} = \max_\omega \big( E_{d \sim p_d}[\log(D_\omega(d))] + E_{x \sim p_x}[\log(1 - D_\omega(E_\theta(x)))] \big). \quad (4)$$

Here $p_x$ is the data distribution and $p(d) = \mathcal{N}(d; 0, I)$ is the prior multivariate Gaussian distribution.

In the *classification phase*, classifiers ($C_\phi$) use the latent code ($z = E_\theta(x)$) as input and minimise standard class entropy loss (using predicted values and target vector containing the labels for three tasks) and error is back-propagated

through the network to update $E_\theta$. The encoder/generator network ($E_\theta$) is updated by optimising the following objective function:

$$L_{enc} = \min_\theta \big( E_{x \sim p_x}[\log(1 - D_\omega(E_\theta(x)))]$$
$$+ E_{x \sim p_x}[\beta \mathcal{L}_{AE}(x, D_\delta(E_\theta(x)))] + \qquad (5)$$
$$E_{x,y \sim p_{X,Y}}[\mathcal{L}_c(E_\theta(x), y; C_\phi)]\big).$$

Unlike the discriminator, the encoder/generator is updated in all three phases. The first term in Equation (5) is updated in the regularisation phase, and the second term in the reconstruction phase and the third term is updated in the classification phase. Also, these three phases run in serial: the reconstruction phase runs first followed by the regularisation and classification phase. In this way, the latent code generation, which is an unsupervised process, gets influenced by the supervised classification task and thus facilitates semi-supervised learning. Note here that when using additional auxiliary data with no labels for emotion, loss functions for gender and speaker are only calculated to update the encoder.

## 4 EXPERIMENTAL SETUP

### 4.1 Datasets

To evaluate the performance of our proposed model, we use two different datasets: IEMOCAP and MSP-IMPROV, which are commonly used for emotion classification research [82], [83]. Both datasets have similar labelling schemes and were collected to simulate naturalistic dyadic interactions between actors. In order to use additional data for the gender and speaker recognition auxiliary tasks, we use Librispeech [84], which is a corpus of read English speech, suitable for training and evaluating speech recognition and speaker identification systems. Below we briefly describe these datasets.

#### 4.1.1 IEMOCAP

This database contains 12 hours of audiovisual data including audio, video, facial motion information, and textual transcriptions [38]. The recordings were collected from 10 professional actors, including five males and five females, during dyadic interactions. This allowed actors to perform spontaneous emotion in contrast to reading text with prototypical emotions [85]. Each interaction is around five minutes long and segmented into smaller utterances of sentences. For categorical labels, each sentence is annotated by three annotators and the participant. Finally, an utterance is given a label if at least three annotators assigned the same label. Overall, this corpus contains nine emotions: angry, excited, happy, sad, neutral, disgust, frustrated, fearful, and surprised. For dimensional annotation, two annotators and the participant were asked to label activation and valence on a scale of 1 to 5. Similarly to prior studies [86], we used utterances of four categorical emotions including angry, happy, neutral, and sad in this study by merging "happy" and "excited" as one emotion class "happy". The final dataset includes 5531 utterances (1103 angry, 1708 neutral, 1084 sad, and 1636 happy).

#### 4.1.2 MSP-IMPROV

The MSP-IMPROV dataset is a multimodal emotional database recorded from 12 actors performing dyadic interactions [39]. The utterances are grouped into six sessions and each session has one male and one female actor similar to IEMOCAP [38]. The scenarios were carefully designed to promote naturalness, while maintaining control over lexical and emotional contents. The emotional labels were collected through perceptual evaluations using crowdsourcing [87]. The utterances in this corpus are annotated on four categorical emotions: angry, happy, neutral, sad. To be consistent with previous studies [15], [88], we use all utterances with four emotions: anger (792), happy (2644), sad (885), and neutral (3477). For dimensional annotation (i.e., activation, and valence), similar to IEMOCAP, these utterances are also annotated on a scale of 1 to 5.

#### 4.1.3 LibriSpeech

The LibriSpeech dataset [84] is derived from audiobooks and it contains 1000 hours of English read speech from 2484 speakers. This corpus is commonly used for speech recognition and speaker identification problems [89], [90]. The training portion of this corpus is split into three subsets, with an approximate recording time of 100, 360 and 500 hours. We used the subset that contains 360 hours of recordings. These recordings span over 961 speakers. Our selection is motivated by the fact it is obviously larger than the 100-hour subset, which spans over only 251 speakers. Also, it offers higher recording quality compared to the 500-hour subset [84]. From this subset, we randomly select 600 speakers (the rationale for choosing the number of speakers is discussed in Section 6.2).

### 4.2 Speech Preprocessing

We have represented the audio utterances in the form of spectrograms, which is a popular 2D representation widely used for speech emotion recognition [81], [91], [92]. The spectrograms were computed using a short-time Fourier transform (STFT) with an overlapping Hamming window of size 25 ms with a 10 ms shift. The height of the spectrogram is 128, which represents the frequency range 0–8 kHz. Due to the varying lengths of the audio samples, the spectrograms vary in width, which poses a problem for the batch processing of the model training. To compensate for this, a context window of 256 frame is applied to create fixed width segments following the procedure used in [51], [93]. Each segmented spectrogram was assigned the emotion label of the corresponding utterance. It is pointed out by previous research that removing silence pauses provides better SER results using deep learning [94], [95]. One of the reasons is that silence adds no speech information to the training data, especially for deep learning models. Nevertheless, we empirically tested that removing silence pauses offers slightly better performance than retaining them. In our experiments, we removed silence pauses from the utterances. We trained all models using segmented spectrograms. In order to calculate the utterance level prediction during the testing phase, posterior probabilities of segments of spectrograms for given utterances were averaged. This is a well known strategy used in SER [93], [96] and also in studies on speaker identification [97].

## 4.3 Model Configuration

Our semi-supervised architecture is illustrated in Figure 1. The encoder part of the autoencoder network consists of three convolutional layers. Each convolutional layer is followed by a pooling layer. These convolutional layers identify emotionally salient regions within the spectrogram and create feature maps. The pooling layer extract highly relevant features by reducing their dimensions. We use max-pooling layer as it offered better performance compared to average pooling during validation. The encoder/generator part encodes the spectrograms into latent code $z$, which has the dimension $16 \times 16 \times 32$. The size of the latent code was determined using the validation set. Here, we use a multivariate Gaussian distribution ($p(d) = \mathcal{N}(d; 0, I)$) with zero mean and unit standard deviation as prior distribution $p(d)$ that we impose on the latent codes $z$ in the regularisation stage. It helps the AAEs to disentangle important attributes of the input data and makes it suitable for speech emotion classification [98]. In SER, using $\mathcal{N}(d; 0, I)$ as prior helps the autoencoder networks to learn the distribution of emotional structures compared to standard autoencoders as validated with variational autoencoders (VAEs) [86] and AAEs [98], [99].

The model was trained with the batch size of 32, where Stochastic Gradient Descent (SGD) with learning rate of 0.0001 was used as optimiser. During validation, accuracy was computed at the end of each epoch. If the accuracy of the model did not improve on the validation set after 5 epochs, we restored the model to best epoch and learning rate was halved. This process continued until the learning rate reached below 0.00001. We applied batch normalisation [100] after each convolutional layer to achieve a stable distribution of activation values throughout the training. The batch normalisation layer was used before the nonlinearity layer. We used a rectified linear unit (ReLU) as non-linear activation function type as it gives us better performance compared to leaky ReLU and hyperbolic tangent during validation. The decoder block has the same structure as the encoder/generator except that the convolutional layers are replaced with transposed convolution layers.

The latent code $z$ was fed to the classifiers, which has four components: (1) convolutional layer, (2) max-pooling layer, (3) dense layers, and (4) softmax layer. We used one convolutional layer followed by max-pooling in each classifier to capture features related to the classification tasks. After each max-pooling layer, we used dense layers followed by a softmax layer to provide prediction. We used two dense layers in each classifier and used a dropout layer, with a dropout rate of 0.3, between them to avoid overfitting. The discriminator of the AAE had a similar architecture to the classifiers, which consists of one convolutional layer followed by a max-pooling layer, and two dense layers followed by a softmax layer.

We performed the step-by-step training of all models. We randomly initialised the model and trained first on Librispeech dataset for speaker and gender classification only. The weights learnt in this stage were used to initialise the autoencoder network when emotional data is fed to the model.

## 5 EXPERIMENTS AND EVALUATIONS

In this study, we evaluated the performance of the proposed framework using 10-fold cross-validation and leave-one-speaker-out validation to compare with multiple studies. For 10-fold cross-validation, we followed the strategies used in [88]. We created the ten folds based on speaker ID so that each fold has all speakers. This allows us to use speaker identification as secondary task. In each step of the validation, one fold was used as validation set for parameter selection, eight folds were used for training, and the remaining fold was used for testing (same as used in [88]).

The 10-fold cross-validation scheme does not allow for speaker-independent testing as such since each fold has data from all speakers. To perform speaker-independent testing we used the leave-one-speaker-out cross-validation scheme commonly used in the literature [30], [86]. This ensures that the speakers are independent in each fold. However, speaker-independent testing limits the use of speaker identification as secondary task. Therefore, we performed speaker verification in this scheme. We consider the d-vector framework [101] for speaker verification, which uses the output of the last hidden layer as speaker representation. During training, the speaker classifier is trained to classify speakers on training data and the evaluation phase involves the extraction of a d-vector from the test utterance using the trained speaker classifier. Then cosine distance is computed between the d-vectors of the test and the claimed speaker. The standard test set of the Librispeech data [84] consisting of 40 speakers along with the test sets of IEMOCAP or MSP-IMPROV were used for enrolment and verification purposes. Equal Error Rate (EER) is used as a measure of performance in the speaker verification system. For SER, we used weighted accuracy (WA) and unweighted accuracy (UA) as comparison measures due to their widely accepted use in studies on speech emotion recognition. For each model used in this work, we repeated the evaluation ten times and calculated the mean and standard deviation of WA and UA.

For evaluations, we considered two types of emotional labels: categorical and dimensional. For categorical emotions, we used four emotions including anger, sad, happy, and neutral. For dimensional emotions, we used two different configurations. Firstly, we manually clustered continuous values of dimensional emotions into three levels. We interpreted activation as low, medium, high, and valence as negative, neutral, positive, as used in [30], [102]. Both, the IEMOCAP and MSP-IMPROV databases are annotated for activation and valence using integral values in the range 1 to 5. Table 2 shows the range of three clusters for activation and valence for both datasets, which has also been adopted from [30], [102].

Table 2: Three Levels of Mapping Rules for IEMOCAP and MSP-IMPROV.

| Corpus | Low/Negative | Medium/Neutral | High/Positive |
|---|---|---|---|
| IEMOCAP | [1,2] | (2,3.5] | (3.5,5] |
| MSP-IMPROV | [1,2.5] | (2.5,3.5] | (3.5,5] |

Secondly, we considered dimensional emotion representation using the valence-activation space, which combines
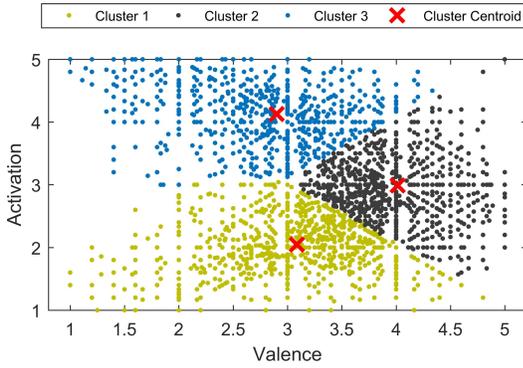
Figure 2: Three clusters of the MSP-IMPROV data in the valence-activation space.

the information of both activation and valence. Previous work has shown that the combination of these two dimensions provides richer emotional information in contrast to using valence and activation separately [103]. To be consistent with previous studies [30], [104], [105], we validated our model on the joint classification of the valence-activation space by building three and five clusters. The cluster midpoints in the valence-activation space were determined by applying $K$-means clustering on the dimensional annotation values of the respective datasets (i.e., IEMOCAP and MSP-IMPROV). A label was assigned to each utterance by choosing the cluster label that minimised the Euclidean distance between the utterance and the cluster centroid. This is highlighted in Figure 2.

### 5.1 Benchmarking Results

We start our evaluation by benchmarking the performance for multi-task over single-task. We implemented a single-task learning (STL) version of our proposed model for a fair comparison. Also, to expand our comparison, we implemented a supervised convolutional neural network (CNN) and an autoencoder (AE) based semi-supervised framework, both using single task learning. We compared the performances of these models for both categorical and dimensional emotions.

As mentioned, while investigating categorical emotion, we use it as the primary task and speaker verification and gender identification as the auxiliary tasks. Then, we use dimensional emotions as the primary task in two ways: grouping continuous values into three levels; and clustering valence-arousal space into three different groups as discussed before. For both ways, the secondary tasks are gender classification and speaker verification.

We trained the single task implementation of our model and the autoencoder in a semi-supervised way for the emotion recognition task. The overall loss function was optimised by tuning the values of $\alpha$ and $\beta$ in Equation (1) and (2) to maximise the system performance and minimise the reconstruction loss, $\mathcal{L}_{\text{AE}}$. Here, Equation (2) contains only first term ($\mathcal{L}_c = \beta * \mathcal{L}_{\text{E}}$) as we are only using emotion classifier. We evaluated the model for different values of $\alpha$ and and $\beta$ ranging from $0.1$ to $1.0$ on the validation set to select the best value. For the IEMOCAP data, we achieve

the best performance for the AE using $\alpha = 0.3$ and $\beta = 0.7$, and for our AAE based models, we achieve the best results for $\alpha = 0.4$ and $\beta = 0.6$.

For two configurations of dimensional emotions, we also identify $\alpha$ and $\beta$ using the validation set and use it for the test set. We achieve the best performance for the AE using $\alpha = 0.6$ and $\beta = 0.4$, and for our AAE based models, we achieve the best results using $\alpha = 0.3$ and $\beta = 0.7$.

The comparisons of results are all summarised in Fig. 3. We observe that, our proposed MTL framework performs better than the STL implementation of our model—the supervised CNN, and the autoencoder. We note this for both categorical and dimensional classification of emotion and for both the IEMOCAP and MSP-IMPROV datasets.

### 5.2 Comparison with Previous Studies

To further extend our comparison scope, in this section, we include results published in recent studies. Note that for IEMOCAP and MSP-IMPROV, there are no standardised training and testing splits to evaluate the results. However, we observe that most of the related studies have used either a 10-fold or a leave-one-speaker-out validation strategy. We therefore implement these schemes and present the comparison results in Table 3. These are, however, accordingly, to be interpreted with the necessary care and merely serve as indication.

Table 3: Comparison of results (UA %) of our proposed method with those of recent studies using categorical emotions.

| Model | IEMOCAP | MSP-IMPROV |
|---|---|---|
| **10-fold cross validation results** | | |
| ProgNet (Transfer Learning) [88] | 65.7 $\pm$1.8 | 60.5 $\pm$2.1 |
| CNN (Muli-task implementation) | 65.6 $\pm$2.0 | 59.5 $\pm$2.4 |
| Semi-supervised AE (Muli-task implementation) | 66.4 $\pm$1.6 | 60.2 $\pm$2.3 |
| Semi-supervised AAE (Proposed) | **68.8$\pm$1.2** | **63.6$\pm$1.7** |
| **leave-one-speaker-out** | | |
| DBN (Multi-task) [30] | 62.2 | – |
| CNN (Multi-task) [106] | 59.54 | – |
| Semi-supervised AAE (proposed) | **66.7$\pm$1.4** | **60.3$\pm$1.1** |
| CVAE-LSTM (Single Task) [20] | 62.8 | – |
| CNN (Single Task) [107] | 64.2 | – |
| Proposed (Single Task) | 64.5$\pm$1.5 | 58.1$\pm$1.7 |

For 10-fold cross validation, we followed the evaluation scheme used in [88]. In [88], authors used progressive neural network and transfer learning (TL) to transfer knowledge from gender and speaker identification to improve the SER performance. Compared to this study, we are achieving better results by exploiting speaker and gender identification as auxiliary tasks within our multi-task learning framework. This shows that, transferring knowledge using auxiliary tasks in MTL can provide more useful information to improve SER performance. In Table 3, we also report the performance of a CNN and an AE when implemented for multi-task learning. We implemented a multitask CNN with two convolutional layers shared with the classification networks following the technique used in [42], [51]. The classification networks consisted of one convolutional layer, two dense layers, and one softmax layer. The AE model had a similar architecture (i.e., hidden units, layers, and model parameters) as our AAE based model—just without the discriminator.

In order to evaluate the proposed model for speaker-independent SER, we used leave-one-speaker-out training
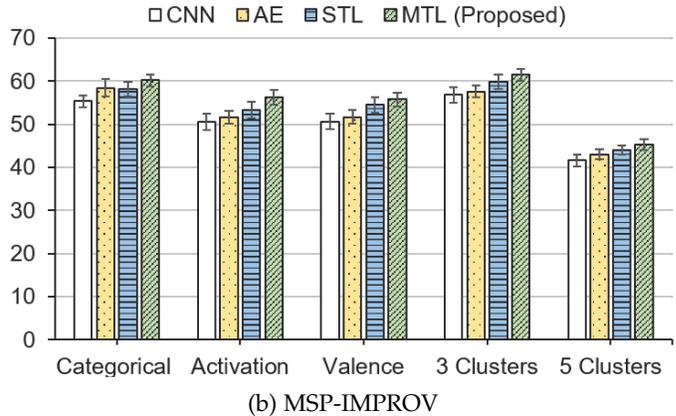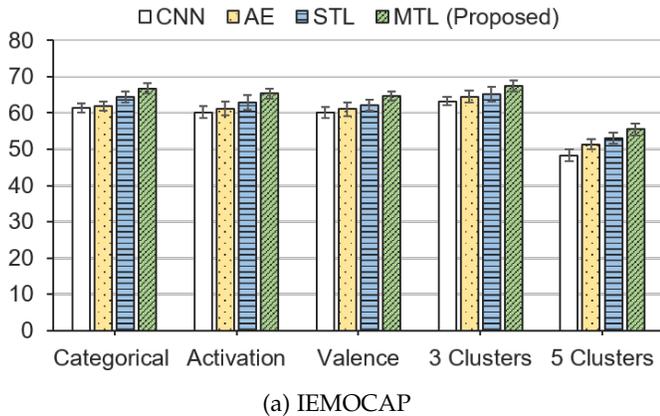
Figure 3: Benchmarking results of the proposed multi-task model (MTL) against a single task implementation of the same model (STL), single task implementation by CNN, and a single-task semi-supervised implementation of an autoencoder (AE) using leave-one-speaker-out scheme

.

with five-fold cross-validation. As speaker independent scheme limits speaker identification as an auxiliary task, therefore, we performed speaker verification. We compare our results with recent studies [20], [30], [107] using speaker-independent SER. In [30], the authors used a multi-task DBN for SER and showed the improved results compared to STL. In [106] the authors used multi-task CNN and utilised additional unlabelled data in an unsupervised way to improve SER performance. Similar to these studies, we also consider MTL framework and achieved better performance than their approaches as reported in Table 3. The authors in [30] and [106] utilised other emotional attributes as auxiliary tasks which limit the use of additional data.

While using the leave-one-speaker-out training with five-fold cross-validation, our proposed model in STL setting also provides better results compared with the other recent studies [20], [107] using STL. Note that these studies did not report any result for the MSP-IMPROV dataset. Also note that, due to the difference in the activation and valence classification strategies, we could not present the results of [51] in Table 3.

### 5.3 Cross-Corpus Results

To verify the generalisability of the proposed model, we also perform a cross-corpus analysis. In this scenario, we trained models using IEMOCAP, and testing is performed on the MSP-IMPROV set. We selected IEMOCAP as training data since it is more balanced and also for good comparison with recent studies, as these studies used a similar scheme [106], [108], [109]. We used 30 % of the MSP-IMPROV data for parameter selection and 70 % as testing data. Here, we used gender classification and speaker verification as an auxiliary task, as speakers in both datasets are different.

We compared our results with other studies on cross-corpus SER. For example, Neumann et al. [106] utilised the representations learnt by autoencoder from unlabelled data fed into a CNN-based classifier. They used the Librispeech and Tedlium (release 2) [110] datasets as unlabelled data, and were able to improve the performance for cross-corpus SER. Our proposed model provides better results compared to this study by using additional data for auxiliary task.

In [108], the authors used Cycle consistent adversarial networks, i.e., the (CycleGAN)-based method to transfer feature vectors extracted from a large unlabelled speech corpus into synthetic features representing the given target emotions. They used Tedlium (release 2) as unlabelled data to generate synthetic data and used this data to augment the classifier. Similarly, Sahu et al. [109] applied generated samples by GANs as additional data to train the classifiers for cross-corpus SER. Both of these studies [108], [109] used additional data to augment the classifiers for cross-corpus SER. However, our approach is different as we are using additional data for auxiliary tasks and achieving similar results without augmenting the system with synthetic data. We compare our results with those of these studies and the comparisons are presented in Table 4. The results show that we achieve competitive accuracy attesting the generalisation ability of the proposed model.

Table 4: Cross-corpus evaluation results for emotion recognition.

| Model | UA (%) |
| --- | --- |
| Attentive CNN [106] | 45.76 |
| Conditional-GAN [109] | 45.40 |
| CycleGAN-DNN [108] | 46.52±0.43 |
| Proposed | 46.41±0.32 |

## 6 ANALYSIS AND DISCUSSION

The experimental results clearly show that the proposed semi-supervised multi-task framework offers an improved performance in speech emotion recognition compared to previous studies. In this section, we focus on three aspects of our proposed model: (1) we elaborate on the impact of a secondary task on improving the performance of the primary task; (2) we quantify the impact of using additional data; and (3) we quantify the impact of tuning the trade-off parameters. All the results in this section are computed using speaker-independent evaluation.

### 6.1 Impact of Secondary Tasks on Primary Task

We consider four categorical emotions from both the IEMO-CAP and MSP-IMPROV datasets as described in Section
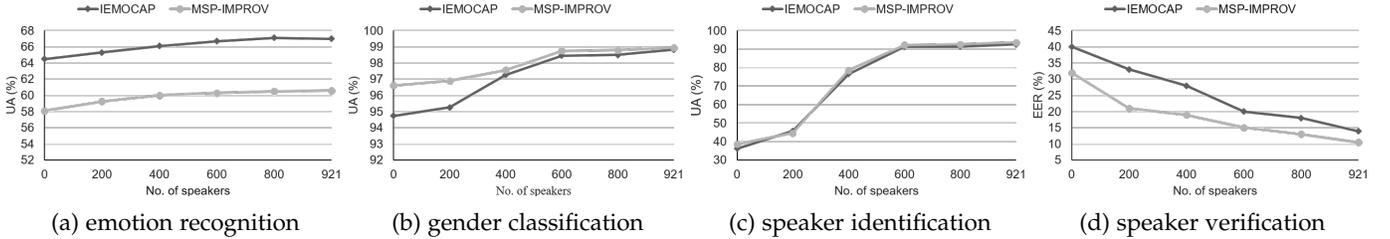
Figure 4: Impact of additional data injection on gender, speaker and emotion recognition.

Table 5: Impact of auxiliary tasks on categorical emotions.

| Secondary Tasks | Primary task: categorical emotion | | | |
| | IEMOCAP | | MSP-IMPROV | |
| | WA | UA | WA | UA |
|---|---|---|---|---|
| Gender | 67.5±1.5 | 66.1±1.7 | 60.1±1.2 | 59.3±1.1 |
| Speaker | 67.2±1.3 | 65.9±1.6 | 60.5±1.2 | 59.1±1.5 |
| Both | **68.5±1.2** | **66.7±1.4** | **62.5±1.4** | **60.2±1.2** |

Table 6: Impact of auxiliary tasks on dimensional emotions.

| Secondary Tasks | Primary Tasks: dimensional emotion | | | |
| | Individual levels of activation and valence | | | |
| | IEMOCAP | | MSP-IMPROV | |
| | Activation | Valence | Activation | Valence |
|---|---|---|---|---|
| Gender | 62.6±1.3 | 61.4±1.2 | 53.7±2.1 | 53.4 ±1.8 |
| Speaker | 63.7 ±1.2 | 60.8±0.8 | 52.6 ±2.5 | 52.6±2.3 |
| Both | **64.5±1.5** | **62.2±1.0** | **54.6±1.4** | **55.4±1.6** |
| Secondary Tasks | Joint activation-valence space | | | |
| | 3 Clusters | 5 Clusters | 3 Clusters | 5 Clusters |
| Gender | 65.7±1.3 | 53.2 ±0.8 | 60.1±1.5 | 43.6±1.5 |
| Speaker | 64.8 ±1.5 | 54.2 ±1.2 | 61.4±1.0 | 43.1±1.2 |
| Both | **65.1±0.9** | **55.1±1.4** | **62.1±1.8** | **44.3±1.3** |

4.1. We also use dimensional emotions in our experiments. In Table 5 and Table 6, we present the results of using the auxiliary tasks separately and jointly to improve the performance of the primary task for categorical and dimensional emotions, respectively. We observe that, while using the auxiliary tasks individually, our model offers similar performance improvement. However, when using the auxiliary tasks jointly, our model offers the highest accuracy for the primary task across both categorical and dimensional emotion representations. Intuitively, jointly learning a representation for emotions with speaker and gender helps to uncover the common high-level discriminative representations, which leads to the performance improvements in the SER system.

## 6.2   Impact of using Additional Data

For the auxiliary tasks of speaker and gender recognition, we use additional data that is not labelled for emotion and show that, when the MTL model is trained with additional data from auxiliary tasks, the performance on the emotion recognition task for both datasets. To further show how performance improves while increasing the amount of data, we trained our model by varying the amount of data for auxiliary tasks. Note that we use the LibriSpeech dataset to introduce additional speakers so, in order to increase the amount of data, we increase the number of speakers.

Fig. 4a shows the effect of varying the amount of additional data on the UA (%) of categorical emotion classification using both datasets. We observe that up to 600 speakers, the performance improvement is quite strong, however, beyond that we observe a plateauing effect. This

is an important observation as it can guide researchers to select a possible operating point when using our suggested method.

To get some further insight into the above improvement, we plot the improvement in auxiliary tasks with the increase of data. Fig. 4b, Fig. 4c, Fig. 4d summarise the results. Here, the performance is calculated using our model in the single task learning mode. We plot the results as we increase the number of speakers. We notice a similar trend as we observe in Fig. 4a. After 600 speakers, improvement in secondary tasks sees a plateau effect. Here, we also plot the speaker verification EER (%) with the increase of number of speakers by using 20 utterances of speakers for enrolment and remaining for evaluation. The performance of speaker verification also improves with the increase of speaker data. Therefore, summarising Fig. 4, it can be noted that improvement in the auxiliary tasks while adding additional data eventually helps to improve the performance of the primary task, which cements the contribution of this paper in proposing a multi-task semi-supervised framework for SER. Intuitively, through the feed of additional data for the auxiliary tasks, a better representation of the intrinsic properties of speech is achieved, which eventually improves the performance of SER.

## 6.3   Impact of Tuning Trade-off Parameters

In this section, we investigate the impact of the trade-off parameters $\alpha$ (Eq (1)), and $\beta$ (Eq (2)), which are the weights of unsupervised, and supervised primary and secondary tasks, on the performance (UA) of the system.

Fig 5 shows the impact of changing the weights $\alpha$ and $\beta$ on UA (%) for categorical, activation, and valence classification. In the first experiment, we keep the value of $\beta$ fixed at $0.5$. This assigns equal weights to both the primary and secondary tasks, but vary the weights (0.1 to 1.0) for the unsupervised task. It can be seen from Fig. 5 that very low and very high weights of $\alpha$ hurt the performance of the system. However, $\alpha$ with values ranging from $0.4 - 0.6$ gave better results for both datasets. This shows that controlling the weights of the unsupervised task through $\alpha$ can improve the performance of the system, however, a suitable range for $\alpha$ needs to be identified that offers the best performance.

Further, Fig. 5 also illustrates the relationship of $\beta$ and UA (%). To highlight this, we select $\alpha = 1$ and vary the weights of $\beta$ (0.1 to 0.9), which controls the weights for both the primary and secondary tasks classification losses (see Equation (2)). It can be noted from Figs. 5a, 5b, and 5c that a very high value of $\beta$ (i. e., the frameworks essentially become single task) gives poor performance. However, we

(a) categorical emotion classification



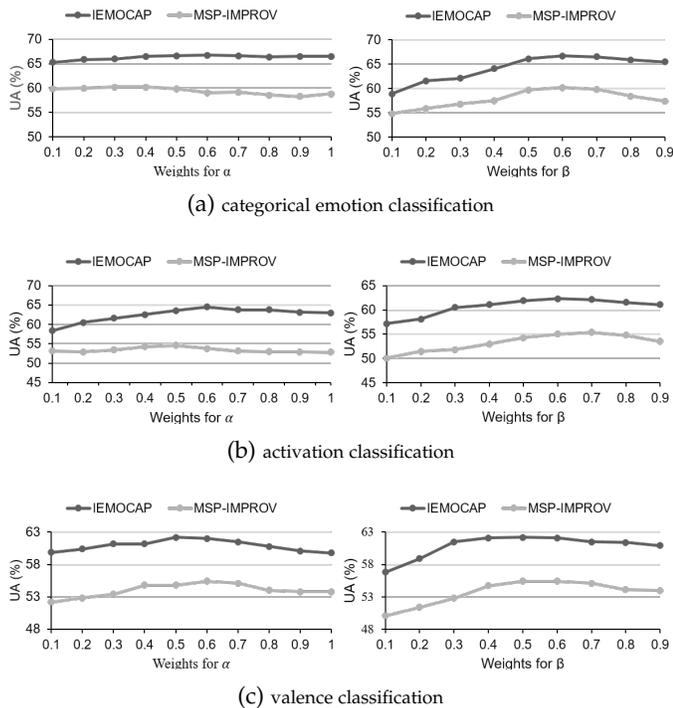(b) activation classification



(c) valence classification

Figure 5: Impact of varying the weights of $\alpha$ and $\beta$ on the performance of categorical (5a), activation (5b), and valence (5c) classification on both, the IEMOCAP and MSP-IMPROV datasets.

also note that too much significance given to auxiliary tasks also diminishes the performance as a very low value of $\beta$ gives poor performance. For both the IEMOCAP and MSP-IMPROV datasets, the system performed better with values of $\beta$ in the range of $0.4 - 0.7$.

## 7 CONCLUSIONS

In contrast with previous studies, this article proposes semi-supervised multi-task learning using adversarial autoencoders for speech emotion recognition (SER). Specifically, we put considerable emphasis on a novel technique of utilising unlabelled data for auxiliary tasks through the proposed multi-task semi-supervised learning model to improve the accuracy of the primary task. We evaluated our proposed model using the popular IEMOCAP and MSP-IMPROV emotion datasets, and demonstrated that it performs notably better than (1) the comparable state-of-the-art studies in SER that use similar methodology and/or implementation strategies; (2) supervised single- and multi-task methods based on CNN, and (3) single- and multi-task semi-supervised autoencoders. We observe this for categorical and dimensional emotion classifications, and cross-corpus SER. Our proposed approach can overcome the challenge of limited data availability of emotion datasets, which is a significant contribution towards developing a robust machine learning model for SER.

Our analysis shows that (1) improvement of the auxiliary tasks through the injection of additional data predominantly drives the improvement of the primary task, (2) a combined effort of auxiliary task is better for improving the accuracy of the primary task, than using them individually, (3) for

the IEMOCAP and MSP-IMPROV datasets, it is possible to reasonably determine an operating point in terms of how much additional data for the auxiliary task is sufficient, (4) it is important to control the weight of loss function of the unsupervised task in the proposed semi-supervised MTL setting to improve the accuracy of SER, and (5) it is important to control the weight of the loss functions of the primary and secondary tasks to achieve the best possible accuracy for SER.

Future work should further focus on the tighter coupling between the generation of data and modelling a richer selection of speaker states and traits simultaneously aiming at 'holistic' speaker analysis [111]. In addition, it appears highly attractive to integrate reinforcement learning into such a framework given a real-life usage such as in a dialogue manager. Likewise, semi-supervised and unsupervised aspects can be benefited by reinforced information.

## REFERENCES

[1] P. Trower, B. Bryant, and M. Argyle, *Social skills and mental health (psychology revivals)*. Routledge, 2013.

[2] F. Burkhardt, J. Ajmera, R. Englert, J. Stegmann, and W. Burleson, "Detecting anger in automated voice portal dialogs," in *Ninth International Conference on Spoken Language Processing*, 2006.

[3] Z. Huang, J. Epps, and D. Joachim, "Speech landmark bigrams for depression detection from naturalistic smartphone speech," in *ICASSP 2019-2019 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 2019, pp. 5856–5860.

[4] R. Rana, S. Latif, R. Gururajan, A. Gray, G. Mackenzie, G. Humphris, and J. Dunn, "Automated screening for distress: A perspective for the future," *European journal of cancer care*, p. e13033, 2019.

[5] K.-Y. Huang, C.-H. Wu, M.-H. Su, and Y.-T. Kuo, "Detecting unipolar and bipolar depressive disorders from elicited speech responses using latent affective structure model," *IEEE Transactions on Affective Computing*, 2018.

[6] M. Merler, K.-N. C. Mac, D. Joshi, Q.-B. Nguyen, S. Hammer, J. Kent, J. Xiong, M. N. Do, J. R. Smith, and R. Feris, "Automatic curation of sports highlights using multimodal excitement features," *IEEE Transactions on Multimedia*, 2018.

[7] H.-J. Vögel, C. Süß, T. Hubregtsen, E. André, B. Schuller, J. Härri, J. Conradt, A. Adi, A. Zadorojniy, J. Terken *et al.*, "Emotion-awareness for intelligent vehicle assistants: A research agenda," in *2018 IEEE/ACM 1st International Workshop on Software Engineering for AI in Autonomous Systems (SEFAIAS)*. IEEE, 2018, pp. 11–15.

[8] L. S. Roberts, "A forensic phonetic study of the vocal responses of individuals in distress," Ph.D. dissertation, University of York, 2012.

[9] N. Ding, V. Sethu, J. Epps, and E. Ambikairajah, "Speaker variability in emotion recognition-an adaptation based approach," in *2012 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 2012, pp. 5101–5104.

[10] T. Vogt and E. André, "Improving automatic emotion recognition from speech via gender differentiation," in *Proc. Language Resources and Evaluation Conference (LREC 2006), Genoa*, 2006.

[11] A. Mill, J. Allik, A. Realo, and R. Valk, "Age-related differences in emotion recognition ability: A cross-sectional study." *Emotion*, vol. 9, no. 5, p. 619, 2009.

[12] S. Latif, A. Qayyum, M. Usman, and J. Qadir, "Cross lingual speech emotion recognition: Urdu vs. western languages," in *2018 International Conference on Frontiers of Information Technology (FIT)*. IEEE, 2018, pp. 88–93.

[13] P. Laukka, D. Neiberg, and H. A. Elfenbein, "Evidence for cultural dialects in vocal emotion expression: Acoustic classification within and across five nations." *Emotion*, vol. 14, no. 3, p. 445, 2014.

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

[15] S. Latif, R. Rana, S. Khalifa, R. Jurdak, and J. Epps, "Direct modelling of speech emotion from raw speech," *arXiv preprint arXiv:1904.03833*, 2019.

[16] S. E. Kahou, C. Pal, X. Bouthillier, P. Froumenty, Ç. Gülçehre, R. Memisevic, P. Vincent, A. Courville, Y. Bengio, R. C. Ferrari *et al.*, "Combining modality specific deep neural networks for emotion recognition in video," in *Proceedings of the 15th ACM on International conference on multimodal interaction*. ACM, 2013, pp. 543–550.

[17] S. Latif, R. Rana, S. Younis, J. Qadir, and J. Epps, "Transfer learning for improving speech emotion classification accuracy," in *Proc. Interspeech 2018*, 2018, pp. 257–261. [Online]. Available: http://dx.doi.org/10.21437/Interspeech.2018-1625

[18] Q. Mao, M. Dong, Z. Huang, and Y. Zhan, "Learning salient features for speech emotion recognition using convolutional neural networks," *IEEE transactions on multimedia*, vol. 16, no. 8, pp. 2203–2213, 2014.

[19] M. Neumann and N. T. Vu, "Attentive convolutional neural network based speech emotion recognition: A study on the impact of input features, signal length, and acted speech," in *Proc. Interspeech 2017*, 2017, pp. 1263–1267. [Online]. Available: http://dx.doi.org/10.21437/Interspeech.2017-917

[20] S. Latif, R. Rana, J. Qadir, and J. Epps, "Variational autoencoders for learning latent representations of speech emotion: A preliminary study," in *Proc. Interspeech 2018*, 2018, pp. 3107–3111. [Online]. Available: http://dx.doi.org/10.21437/Interspeech.2018-1568

[21] S. Mirsamadi, E. Barsoum, and C. Zhang, "Automatic speech emotion recognition using recurrent neural networks with local attention," in *2017 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 2017, pp. 2227–2231.

[22] M. Wöllmer, A. Metallinou, N. Katsamanis, B. Schuller, and S. Narayanan, "Analyzing the memory of blstm neural networks for enhanced emotion classification in dyadic spoken interactions," in *2012 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 2012, pp. 4157–4160.

[23] Z. Zhang, B. Wu, and B. Schuller, "Attention-augmented end-to-end multi-task learning for emotion prediction from speech," in *ICASSP 2019-2019 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 2019, pp. 6705–6709.

[24] R. Caruana, "Multitask learning," *Machine learning*, vol. 28, no. 1, pp. 41–75, 1997.

[25] X. Li, Y.-Y. Wang, and G. Tur, "Multi-task learning for spoken language understanding with shared slots," in *Twelfth Annual Conference of the International Speech Communication Association*, 2011.

[26] R. Collobert and J. Weston, "A unified architecture for natural language processing: Deep neural networks with multitask learning," in *Proceedings of the 25th international conference on Machine learning*. ACM, 2008, pp. 160–167.

[27] J. Kim, G. Englebienne, K. P. Truong, and V. Evers, "Towards speech emotion recognition" in the wild" using aggregated corpora and deep multi-task learning," *arXiv preprint arXiv:1708.03920*, 2017.

[28] S. Parthasarathy and C. Busso, "Jointly predicting arousal, valence and dominance with multi-task learning," *INTERSPEECH, Stockholm, Sweden*, 2017.

[29] R. Lotfian and C. Busso, "Predicting categorical emotions by jointly learning primary and secondary emotions through multitask learning," in *Proc. Interspeech 2018*, 2018, pp. 951–955. [Online]. Available: http://dx.doi.org/10.21437/Interspeech.2018-2464

[30] R. Xia and Y. Liu, "A multi-task learning framework for emotion recognition using 2d continuous space," *IEEE Transactions on Affective Computing*, no. 1, pp. 3–14, 2017.

[31] E. Hudlicka, "Computational modeling of cognition–emotion interactions: Theoretical and practical relevance for behavioral healthcare," in *Emotions and Affect in Human Factors and Human-Computer Interaction*. Elsevier, 2017, pp. 383–436.

[32] P. Ekman, "An argument for basic emotions," *Cognition & emotion*, vol. 6, no. 3-4, pp. 169–200, 1992.

[33] J. Posner, J. A. Russell, and B. S. Peterson, "The circumplex model of affect: An integrative approach to affective neuroscience, cognitive development, and psychopathology," *Development and psychopathology*, vol. 17, no. 3, pp. 715–734, 2005.

[34] X. Liu, J. Van De Weijer, and A. D. Bagdanov, "Exploiting unlabeled data in cnns by self-supervised learning to rank," *IEEE transactions on pattern analysis and machine intelligence*, 2019.

[35] X. Liu, J. van de Weijer, and A. D. Bagdanov, "Leveraging unlabeled data for crowd counting by learning to rank," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2018, pp. 7661–7669.

[36] S. Sahu, R. Gupta, G. Sivaraman, W. AbdAlmageed, and C. Espy-Wilson, "Adversarial auto-encoders for speech based emotion recognition," *arXiv preprint arXiv:1806.02146*, 2018.

[37] A. Makhzani, J. Shlens, N. Jaitly, I. Goodfellow, and B. Frey, "Adversarial autoencoders," *ICLR 2016 Workshop*, 2015.

[38] 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.

[39] C. Busso, S. Parthasarathy, A. Burmania, M. AbdelWahab, N. Sadoughi, and E. M. Provost, "Msp-improv: An acted corpus of dyadic interactions to study emotion perception," *IEEE Transactions on Affective Computing*, no. 1, pp. 67–80, 2017.

[40] S. Ben-David and R. Schuller, "Exploiting task relatedness for multiple task learning," in *Learning Theory and Kernel Machines*. Springer, 2003, pp. 567–580.

[41] J. Baxter, "A model of inductive bias learning," *Journal of Artificial Intelligence Research*, vol. 12, pp. 149–198, 2000.

[42] S. Li, Z.-Q. Liu, and A. B. Chan, "Heterogeneous multi-task learning for human pose estimation with deep convolutional neural network," in *Proceedings of the IEEE conference on computer vision and pattern recognition workshops*, 2014, pp. 482–489.

[43] Z. Zhang, P. Luo, C. C. Loy, and X. Tang, "Facial landmark detection by deep multi-task learning," in *European Conference on Computer Vision*. Springer, 2014, pp. 94–108.

[44] C. Farabet, C. Couprie, L. Najman, and Y. LeCun, "Learning hierarchical features for scene labeling," *IEEE transactions on pattern analysis and machine intelligence*, vol. 35, no. 8, pp. 1915–1929, 2013.

[45] D. Chen, S. Ren, Y. Wei, X. Cao, and J. Sun, "Joint cascade face detection and alignment," in *European Conference on Computer Vision*. Springer, 2014, pp. 109–122.

[46] X. Zhu and D. Ramanan, "Face detection, pose estimation, and landmark localization in the wild," in *Computer Vision and Pattern Recognition (CVPR), 2012 IEEE Conference on*. IEEE, 2012, pp. 2879–2886.

[47] R. Ranjan, V. M. Patel, and R. Chellappa, "Hyperface: A deep multi-task learning framework for face detection, landmark localization, pose estimation, and gender recognition," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 41, no. 1, pp. 121–135, 2019.

[48] D. Le, Z. Aldeneh, and E. M. Provost, "Discretized continuous speech emotion recognition with multi-task deep recurrent neural network," *Interspeech, 2017 (to apear)*, 2017.

[49] F. Ma, W. Gu, W. Zhang, S. Ni, S.-L. Huang, and L. Zhang, "Speech emotion recognition via attention-based dnn from multi-task learning," in *Proceedings of the 16th ACM Conference on Embedded Networked Sensor Systems*. ACM, 2018, pp. 363–364.

[50] F. Eyben, M. Wöllmer, and B. Schuller, "A multitask approach to continuous five-dimensional affect sensing in natural speech," *ACM Transactions on Interactive Intelligent Systems (TiiS)*, vol. 2, no. 1, p. 6, 2012.

[51] J. Chang and S. Scherer, "Learning representations of emotional speech with deep convolutional generative adversarial networks," in *Acoustics, Speech and Signal Processing (ICASSP), 2017 IEEE International Conference on*. IEEE, 2017, pp. 2746–2750.

[52] J. Carletta, S. Ashby, S. Bourban, M. Flynn, M. Guillemot, T. Hain, J. Kadlec, V. Karaiskos, W. Kraaij, M. Kronenthal *et al.*, "The ami meeting corpus: A pre-announcement," in *International workshop on machine learning for multimodal interaction*. Springer, 2005, pp. 28–39.

[53] F. Tao and G. Liu, "Advanced lstm: A study about better time dependency modeling in emotion recognition," in *2018 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 2018, pp. 2906–2910.

[54] B. Zhang, E. M. Provost, and G. Essl, "Cross-corpus acoustic emotion recognition with multi-task learning: Seeking common ground while preserving differences," *IEEE Transactions on Affective Computing*, no. 1, pp. 1–1, 2017.

[55] Z. Zhang, J. Han, J. Deng, X. Xu, F. Ringeval, and B. Schuller, "Leveraging unlabeled data for emotion recognition with en-

hanced collaborative semi-supervised learning," *IEEE Access*, vol. 6, pp. 22 196–22 209, 2018.

[56] J. Deng, X. Xu, Z. Zhang, S. Frühholz, and B. Schuller, "Semisupervised autoencoders for speech emotion recognition," *IEEE/ACM Transactions on Audio, Speech, and Language Processing*, vol. 26, no. 1, pp. 31–43, 2018.

[57] J. Huang, Y. Li, J. Tao, Z. Lian, M. Niu, and J. Yi, "Speech emotion recognition using semi-supervised learning with ladder networks," in *2018 First Asian Conference on Affective Computing and Intelligent Interaction (ACII Asia)*. IEEE, 2018, pp. 1–5.

[58] J.-H. Tao, J. Huang, Y. Li, Z. Lian, and M.-Y. Niu, "Semi-supervised ladder networks for speech emotion recognition," *International Journal of Automation and Computing*, 2019.

[59] S. Parthasarathy and C. Busso, "Semi-supervised speech emotion recognition with ladder networks," *arXiv preprint arXiv:1905.02921*, 2019.

[60] Z. Huang, M. Dong, Q. Mao, and Y. Zhan, "Speech emotion recognition using cnn," in *Proceedings of the 22Nd ACM International Conference on Multimedia*, 2014.

[61] P. Jackson and S. Haq, "Surrey audio-visual expressed emotion(savee) database," *University of Surrey: Guildford, UK*, 2014.

[62] F. Burkhardt, A. Paeschke, M. Rolfes, W. F. Sendlmeier, and B. Weiss, "A database of german emotional speech." in *Interspeech*, vol. 5, 2005, pp. 1517–1520.

[63] I. Engberg and A. Hansen, "Documentation of the danish emotional speech database des," *Aalborg*, 1996.

[64] L. Fu, X. Mao, and L. Chen, "Speaker independent emotion recognition based on svm/hmms fusion system," in *2008 International Conference on Audio, Language and Image Processing*. IEEE, 2008, pp. 61–65.

[65] S. Parthasarathy and C. Busso, "Ladder networks for emotion recognition: Using unsupervised auxiliary tasks to improve predictions of emotional attributes," *Proc. Interspeech 2018*, pp. 3698–3702, 2018.

[66] H. Larochelle, M. Mandel, R. Pascanu, and Y. Bengio, "Learning algorithms for the classification restricted boltzmann machine," *Journal of Machine Learning Research*, vol. 13, no. Mar, pp. 643–669, 2012.

[67] 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.

[68] A. Makhzani, "Unsupervised representation learning with autoencoders," Ph.D. dissertation, 2018.

[69] J. Deng, Z. Zhang, and B. Schuller, "Linked source and target domain subspace feature transfer learning–exemplified by speech emotion recognition," in *2014 22nd International Conference on Pattern Recognition*. IEEE, 2014, pp. 761–766.

[70] A. Krizhevsky, I. Sutskever, and G. E. Hinton, "Imagenet classification with deep convolutional neural networks," in *Advances in neural information processing systems*, 2012, pp. 1097–1105.

[71] Y. LeCun, B. E. Boser, J. S. Denker, D. Henderson, R. E. Howard, W. E. Hubbard, and L. D. Jackel, "Handwritten digit recognition with a back-propagation network," in *Advances in neural information processing systems*, 1990, pp. 396–404.

[72] S. Lawrence, C. L. Giles, A. C. Tsoi, and A. D. Back, "Face recognition: A convolutional neural-network approach," *IEEE transactions on neural networks*, vol. 8, no. 1, pp. 98–113, 1997.

[73] X. Zhang, J. Zhao, and Y. LeCun, "Character-level convolutional networks for text classification," in *Advances in neural information processing systems*, 2015, pp. 649–657.

[74] W. Xiong, L. Wu, F. Alleva, J. Droppo, X. Huang, and A. Stolcke, "The microsoft 2017 conversational speech recognition system," in *2018 IEEE international conference on acoustics, speech and signal processing (ICASSP)*. IEEE, 2018, pp. 5934–5938.

[75] Z. Huang, M. Dong, Q. Mao, and Y. Zhan, "Speech emotion recognition using cnn," in *Proceedings of the 22nd ACM international conference on Multimedia*. ACM, 2014, pp. 801–804.

[76] J. Liu, W. Han, H. Ruan, X. Chen, D. Jiang, and H. Li, "Learning salient features for speech emotion recognition using cnn," in *2018 First Asian Conference on Affective Computing and Intelligent Interaction (ACII Asia)*. IEEE, 2018, pp. 1–5.

[77] M. Neumann and N. T. Vu, "Attentive convolutional neural network based speech emotion recognition: A study on the impact of input features, signal length, and acted speech," *Proc. Interspeech 2017*, pp. 1263–1267, 2017.

[78] G. Trigeorgis, F. Ringeval, R. Brueckner, E. Marchi, M. A. Nicolaou, B. Schuller, and S. Zafeiriou, "Adieu features? end-to-end speech emotion recognition using a deep convolutional recurrent network," in *2016 IEEE international conference on acoustics, speech and signal processing (ICASSP)*. IEEE, 2016, pp. 5200–5204.

[79] W. Lim, D. Jang, and T. Lee, "Speech emotion recognition using convolutional and recurrent neural networks," in *2016 Asia-Pacific Signal and Information Processing Association Annual Summit and Conference (APSIPA)*. IEEE, 2016, pp. 1–4.

[80] J. Zhao, X. Mao, and L. Chen, "Speech emotion recognition using deep 1d & 2d cnn lstm networks," *Biomedical Signal Processing and Control*, vol. 47, pp. 312–323, 2019.

[81] C. Etienne, G. Fidanza, A. Petrovskii, L. Devillers, and B. Schmauch, "Cnn+ lstm architecture for speech emotion recognition with data augmentation," in *Proc. Workshop on Speech, Music and Mind 2018*, 2018, pp. 21–25.

[82] R. Lotfian and C. Busso, "Retrieving categorical emotions using a probabilistic framework to define preference learning samples," in *Interspeech 2016*, 2016, pp. 490–494. [Online]. Available: http://dx.doi.org/10.21437/Interspeech.2016-1052

[83] Y. Kim and E. M. Provost, "Emotion spotting: Discovering regions of evidence in audio-visual emotion expressions," in *Proceedings of the 18th ACM International Conference on Multimodal Interaction*. ACM, 2016, pp. 92–99.

[84] V. Panayotov, G. Chen, D. Povey, and S. Khudanpur, "Librispeech: an asr corpus based on public domain audio books," in *2015 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 2015, pp. 5206–5210.

[85] R. Lotfian and C. Busso, "Building naturalistic emotionally balanced speech corpus by retrieving emotional speech from existing podcast recordings," *IEEE Transactions on Affective Computing*, no. 1, pp. 1–1, 2017.

[86] S. Latif, R. Rana, J. Qadir, and J. Epps, "Variational autoencoders for learning latent representations of speech emotion: A preliminary study," *Proc. Interspeech 2018*, pp. 3107–3111, 2018.

[87] 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, 2016.

[88] 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.

[89] A. Bérard, L. Besacier, A. C. Kocabiyikoglu, and O. Pietquin, "End-to-end automatic speech translation of audiobooks," in *2018 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 2018, pp. 6224–6228.

[90] H. Dubey, A. Sangwan, and J. H. Hansen, "Transfer learning using raw waveform sincnet for robust speaker diarization," in *ICASSP 2019-2019 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 2019, pp. 6296–6300.

[91] N. Cummins, S. Amiriparian, G. Hagerer, A. Batliner, S. Steidl, and B. Schuller, "An image-based deep spectrum feature representation for the recognition of emotional speech," in *Proceedings of the 25th ACM international conference on Multimedia*. ACM, 2017, pp. 478–484.

[92] S. Albanie, A. Nagrani, A. Vedaldi, and A. Zisserman, "Emotion recognition in speech using cross-modal transfer in the wild," *arXiv preprint arXiv:1808.05561*, 2018.

[93] P. Yenigalla, A. Kumar, S. Tripathi, C. Singh, S. Kar, and J. Vepa, "Speech emotion recognition using spectrogram & phoneme embedding." in *Interspeech*, 2018, pp. 3688–3692.

[94] Z.-Q. Wang and I. Tashev, "Learning utterance-level representations for speech emotion and age/gender recognition using deep neural networks," in *2017 IEEE international conference on acoustics, speech and signal processing (ICASSP)*. IEEE, 2017, pp. 5150–5154.

[95] A. Satt, S. Rozenberg, and R. Hoory, "Efficient emotion recognition from speech using deep learning on spectrograms." in *INTERSPEECH*, 2017, pp. 1089–1093.

[96] S. Latif, J. Qadir, and M. Bilal, "Unsupervised adversarial domain adaptation for cross-lingual speech emotion recognition," in *2019 8th International Conference on Affective Computing and Intelligent Interaction (ACII)*. IEEE, 2019, pp. 732–737.

[97] M. Ravanelli and Y. Bengio, "Speaker recognition from raw waveform with sincnet," in *2018 IEEE Spoken Language Technology Workshop (SLT)*. IEEE, 2018, pp. 1021–1028.

[98] S. Parthasarathy, V. Rozgic, M. Sun, and C. Wang, "Improving emotion classification through variational inference of latent variables," in *ICASSP 2019-2019 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 2019, pp. 7410–7414.

[99] S. E. Eskimez, Z. Duan, and W. Heinzelman, "Unsupervised learning approach to feature analysis for automatic speech emotion recognition," in *2018 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 2018, pp. 5099–5103.

[100] S. Ioffe and C. Szegedy, "Batch normalization: Accelerating deep network training by reducing internal covariate shift," *arXiv preprint arXiv:1502.03167*, 2015.

[101] E. Variani, X. Lei, E. McDermott, I. L. Moreno, and J. Gonzalez-Dominguez, "Deep neural networks for small footprint text-dependent speaker verification," in *2014 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 2014, pp. 4052–4056.

[102] M. Wöllmer, A. Metallinou, F. Eyben, B. Schuller, and S. Narayanan, "Context-sensitive multimodal emotion recognition from speech and facial expression using bidirectional lstm modeling," in *Proc. INTERSPEECH 2010, Makuhari, Japan*, 2010, pp. 2362–2365.

[103] C.-C. Lee, C. Busso, S. Lee, and S. S. Narayanan, "Modeling mutual influence of interlocutor emotion states in dyadic spoken interactions," in *Tenth Annual Conference of the International Speech Communication Association*, 2009.

[104] 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, 2013.

[105] A. Metallinou, M. Wollmer, A. Katsamanis, F. Eyben, B. Schuller, and S. Narayanan, "Context-sensitive learning for enhanced audiovisual emotion classification," *IEEE Transactions on Affective Computing*, vol. 3, no. 2, pp. 184–198, 2012.

[106] M. Neumann and N. T. Vu, "Improving speech emotion recognition with unsupervised representation learning on unlabeled speech," in *ICASSP 2019-2019 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 2019, pp. 7390–7394.

[107] X. Ma, Z. Wu, J. Jia, M. Xu, H. Meng, and L. Cai, "Emotion recognition from variable-length speech segments using deep learning on spectrograms," *Proc. Interspeech 2018*, pp. 3683–3687, 2018.

[108] F. Bao, M. Neumann, and N. T. Vu, "Cyclegan-based emotion style transfer as data augmentation for speech emotion recognition," *Manuscript submitted for publication*, pp. 35–37, 2019.

[109] S. Sahu, R. Gupta, and C. Espy-Wilson, "On enhancing speech emotion recognition using generative adversarial networks," *Proc. Interspeech 2018*, pp. 3693–3697, 2018.

[110] A. Rousseau, P. Deléglise, and Y. Esteve, "Enhancing the ted-lium corpus with selected data for language modeling and more ted talks." in *LREC*, 2014, pp. 3935–3939.

[111] Y. Zhang, Y. Liu, F. Weninger, and B. Schuller, "Multi-task deep neural network with shared hidden layers: Breaking down the wall between emotion representations," in *Acoustics, Speech and Signal Processing (ICASSP), 2017 IEEE International Conference on*. IEEE, 2017, pp. 4990–4994.