

DNN-HMM based Speaker Adaptive Emotion Recognition using Proposed Epoch and MFCC Features

Md. Shah Fahad^a, Jainath Yadav^b, Gyadhar Pradhan^a, Akshay Deepak^a

^aDepartment of Computer Science, National Institute of Technology Patna, India

^bDepartment of Computer Science, Central University of South Bihar, Patna, India

Abstract

Speech is produced when time varying vocal tract system is excited with time varying excitation source. Therefore, the information present in a speech such as message, emotion, language, speaker is due to the combined effect of both excitation source and vocal tract system. However, there is very less utilization of excitation source features to recognize emotion. In our earlier work, we have proposed a novel method to extract glottal closure instants (GCIs) known as epochs. In this paper, we have explored epoch features namely instantaneous pitch, phase and strength of epochs for discriminating emotions. We have combined the excitation source features and the well known Male-frequency cepstral coefficient (MFCC) features to develop an emotion recognition system with improved performance. DNN-HMM speaker adaptive models have been developed using MFCC, epoch and combined features. IEMOCAP emotional database has been used to evaluate the models. The average accuracy for emotion recognition system when using MFCC and epoch features separately is 59.25% and 54.52% respectively. The recognition performance improves to 64.2% when MFCC and epoch features are combined.

Keywords: Emotion Recognition, Epoch Features, Deep Neural Network(DNN), Gaussian Mixture Model (GMM), Hidden Markov Model(HMM), Zero Time Windowing (ZTW)

1. Introduction

Automatic emotion recognition from speech signal has fascinated the research community in the recent years due to its applicability in real-life. Human beings use a lot of emotions along with textual messages to convey the intended information. Emotions improve human computer interactions (HCI) system such as interactive movies [1], story telling and E-tutoring applications [2], and, retrieval and indexing of the video/audio files [3]. Emotion recognition system assists to improve the quality of service of call attendants at call centers [4]. Automatic emotion detection could be helpful in the psychological treatment as used in references [[5],[6],[7]]. It can also be useful in the case of surveillance systems [8]. Modern speech-based systems are designed largely using neutral speech. Here, the components of emotions can be used as an add-on to improve the accuracy in practical applications.

Excitation source features are not much exploited to recognize emotions. Observation from the literature reveals that the majority of the previous works used prosodic and system features for emotion recognition using speech [[9, 10]]. The system features MFCCs, Linear Predictive Cepstral Coefficients (LPCCs) and their derivatives reflect the emotion specific information. Prosodic features such as fundamental frequency, duration, energy and intonation are also used for emotion recognition. Combinations of prosodic and system features are also

widely used for emotion recognition. Reference [11] uses supra-segmental features such as energy, F0, formant locations, energy, dynamics of F0 and formant contours for emotion classification. The statistical parameters of F0 like maximum, minimum, and median values, and the slopes of F0 contours have emotion specific information [12]. However, not much work has been done in using excitation source features for emotion recognition.

Reference [9] combined 55 features (24 MFCCs, 25 prosodic and 6 formant frequencies) for recognizing six emotions. Prosodic and spectral features are combined in reference [10] for emotion classification. It is proven from literature that a combination of different complement features improve the accuracy of emotion recognition system. Most of the features are extracted from speech based on the assumption that the speech signal is stationary in the small speech segment. However, the speech features – either source features or system features – vary rapidly in emotional speech because of the rapid changes in the vibration of the vocal cords. In reference [13], emotion recognition model is developed using a combination of epoch and MFCC features. The proposed model used zero frequency filter (ZFF) method for extracting epoch features. The accuracy of epoch detection using ZFF decreases for emotional speech because it requires a priory pitch period to detect epoch location. However, the pitch period of emotional speech varies frequently in an utterance. The emotion recognition model (in reference [13]) was developed using auto-associative neural networks (AANN) and support vector machines (SVM) on IITKGP-SESC database.

Email addresses: shah.cse16@nitp.ac.in (Md. Shah Fahad), jainath@cub.ac.in (Jainath Yadav), gdp@nitp.ac.in (Gyadhar Pradhan), akshayd@nitp.ac.in (Akshay Deepak)

In our earlier work [14], we proposed a robust method to detect epoch locations. In this paper, epoch features namely instantaneous pitch, phase and strength of excitation (SOE) are extracted. These features are explored for different emotions and combined with MFCCs for classifying four emotions. Using this method, a significant increase in the accuracy of emotion recognition model was observed. The average accuracy for emotion recognition system when using MFCC and epoch features separately is 59.25% and 54.52% respectively. This improves to 64.2% when MFCC and epoch features are combined.

The rest of the paper is organized as follows. Section 2 contains the description of speech databases, Sec. 3 describes detection of epoch features and Sec. 4 briefly discusses MFCC and development of emotion recognition models. The results are discussed in Sec. 5. Section 6 concludes the paper.

2. Databases

Our proposed model has been evaluated on IEMOCAP (Interactive emotional dyadic motion capture database)[15] and IITKGP:SEHSC (Indian Institute of Technology Kharagpur: Simulated Emotion Hindi Speech Corpus) [16]. IEMOCAP database is a multi-modal database which contains audio, video, text and gesture information of conversations arranged in dyadic sessions. The database is recorded with ten actors (five male and five female) in five sessions. In each session, there are conversations of two actors, one from each gender, on two subjects. The conversation of one session is approximately five minutes long. The contents of the database are recorded in both scripted and spontaneous scenarios. The total number of utterances in the database are 10,039, where 4,784 utterances are from the spontaneous sessions and 5,225 are from the scripted sessions. The average duration of an utterance is 4.5 seconds while the average word count per utterance is 11.4 words. The duration of the database is about 12 hours. The database is labeled as per the two popular schemes: discrete categorical label (i.e, labeled as happy, anger, neutral and sad) and continuous dimensional label (i.e, valence, activation and dominance). We have only used the audio tracks and the corresponding discrete categorical labels for emotion recognition.

In IITKGP-SESC, fifteen emotionally neutral Hindi text prompts were used for recording the emotion in multiple sessions to capture diversity. In each session, 15 sentences in eight basic emotions are uttered by each artist. Recording was done with the help of SHURE dynamic cardioid microphone C660N at 16 kHz sampling frequency. The Hindi emotional speech database has 10 speakers (five males and five females) and 15 sentences were recorded for eight emotions (Neutral, Happy, Angry, Sad, Disgust, Sarcastic, Surprise and Fear). There are a total of 12000 speech utterances (10 speakers x 15 sentences x 8 emotions x 10 sessions) in the Hindi emotional speech database. There are 1500 articulations for each emotions. The number of syllables and words in the sentences lie in the range of 9-17 and 4-7 respectively.

3. Extraction of Epoch features using Zero time Windowing method

In our method, voiced regions are detected using the phase of zero frequency filtered speech signal [17]. After that, Zero Time Windowing (ZTW) method [18] is applied to get Hilbert envelope of the Numerator Group Delay (HNGD) spectra of each of the voiced segments. The amplitude of the sum of the three prominent peaks is obtained from each spectrum of the HNGD. The resulting output reproduces the instantaneous energy profile of the windowed signal. The spectral energy profile, obtained from HNGD spectrum, shows high energy at the epoch locations because of high SNR (signal to noise ratio) at these locations. Further, the spectral energy profile is normalized using mean smooth filter. The normalized spectral energy profile is then convolved with a Gaussian filter to highlight the peaks. The positive peaks – selected after removing the spurious peaks – are considered as epochs. Next, each of the above step is described in detail.

3.1. Voiced Activity Detection (VAD)

Epochs are present in the voiced regions due to vibration of the vocal cords. Hence, we first divide the speech into voiced and unvoiced regions based on its characteristics. In the present paper, voiced regions are detected [17] using the phase of Zero Frequency Filtered Signal (ZFFS). The ZFFS of a speech utterance is obtained by using zero frequency resonator [19]. The phase of a ZFFS is determined using the Hilbert transformation. Further, the phase-signal is split into frames of size 30 ms with frame shift of 5 ms and each frame is convolved with Hanning-window. The amplitude spectrum of Hanning-windowed frame is computed. Thereafter, the sum of the first 10 harmonics is computed. The decision of voiced and unvoiced regions is taken based on the appropriate threshold of global maxima of the sum of phase harmonics (SPH) because the global maxima of the SPH of voiced regions is significantly higher than unvoiced regions.

Voiced and unvoiced regions of a speech signal are detected by setting the threshold of 0.08 for global maxima of SPH of each and every frame as shown in Fig. 1. Fig. 1(a) shows the the speech signal. The corresponding global maxima of SPH is shown in Fig. 1(b), which is separated as voiced and unvoiced speech in Fig. 1(c) through rectangular waveform. Here, voiced speech is labeled 1(high) and unvoiced speech is labeled 0 (low).

3.2. Sequence of Steps for Epoch extraction

The steps to detect epoch locations are described next.

1. The voiced segment is detected using the phase of zero frequency filtered speech signal [17].
2. The voiced speech signal is differentiated to remove any low frequency bias in the speech signal using the formula

$$y[n] = s[n] - s[n - 1] \quad (1)$$

where:

$y[n]$ is the differentiated signal at n^{th} sample

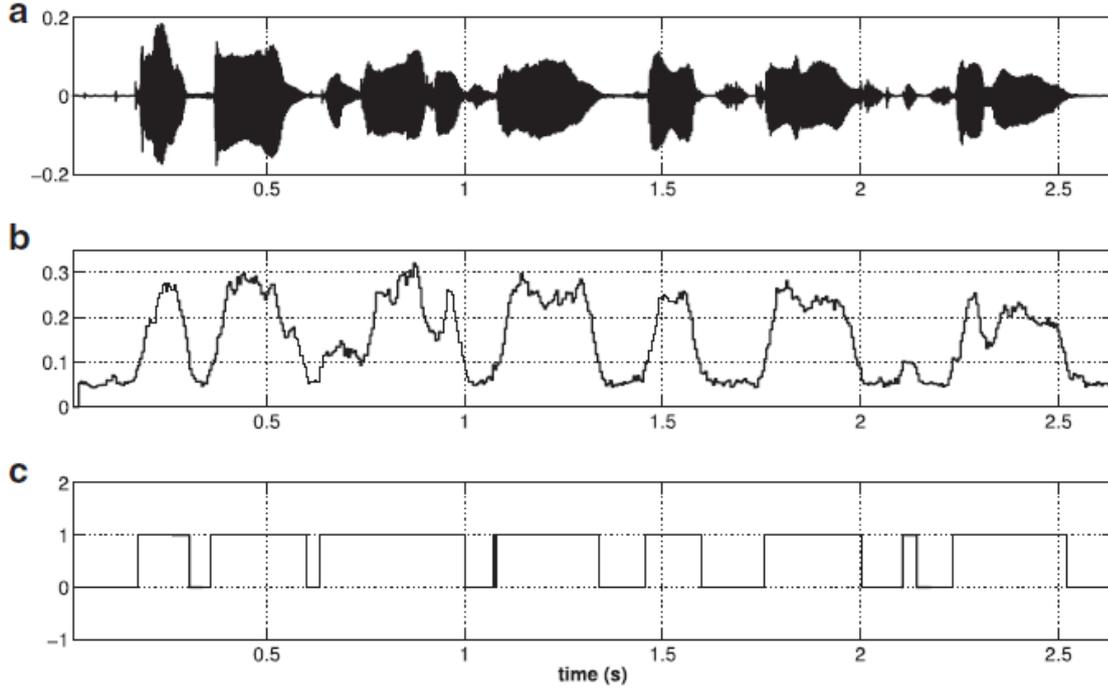


Fig. 1: Detection of voiced and unvoiced regions using the phase of ZFFS. (a) Speech signal. (b) its corresponding global maxima of SPH. (c) unvoiced and voiced regions correspond to low and high amplitude respectively.

$s[n]$ is the actual speech signal at n^{th} sample, and, $s[n - 1]$ is the actual speech signal at $(n - 1)^{\text{th}}$ sample

3. Three milliseconds segments of the differentiated speech signal (resulting in $M = 48$ samples) were taken at each sampling point. These were appended with $N - M$ (2048-48) zeros to obtain sufficient resolution in the frequency domain.
4. The time domain signal is multiplied with the square of window function h_1 (defined below) to achieve the smoothed spectrum by integration in the frequency domain.

$$h_1[n] = \begin{cases} 0 & n = 0 \\ h_1[n] = \frac{1}{4\sin^2(\frac{\pi n}{N})} & n = 1, 2, \dots, N - 1 \end{cases} \quad (2)$$

5. The ripple effect due to truncation is reduced by multiplying the signal of the previous step with the window h_2 , which is defined as:

$$h_2[n] = 4\cos^2\left(\frac{\pi n}{2M}\right), n = 0, 1, 2, \dots, M - 1 \quad (3)$$

The resultant signal $x[n]$ is called windowed signal.

6. To highlight the spectral features, the numerator of group delay of windowed signal, denoted $g[k]$, is computed as:

$$g[k] = X_R[k]Y_R[k] + X_I[k]Y_I[k], k = 0, 1, 2, \dots, N - 1 \quad (4)$$

The resultant signal is known as DNGD signal.

7. Hilbert envelope of the DNGD spectrum is computed to prominently highlight the spectral peaks. The Hilbert envelope $h_e[k]$ of DNGD signal $g[k]$ is computed as:

$$h_e[k] = \sqrt{g^2[k] + g_h^2[k]} \quad (5)$$

where $g_h[k]$ is the Hilbert transformation of the sequence $g[k]$. It is computed as:

$$g_h[k] = IDFT E_h(w) \quad (6)$$

where $E(\omega)$ is the DTFT of the sequence $g(k)$. It is defined as:

$$E_h(\omega) = \begin{cases} -jE(\omega), & 0 < \omega < \pi \\ jE(\omega), & -\pi < \omega < 0 \end{cases} \quad (7)$$

8. The sum of the three most prominent peaks of the HNGD spectrum is determined at each sampling instant. The resultant amplitude shows high SNR around glottal closure. Further, the amplitude contour is smoothed using 5-point mean smoothing filter to eliminate any outliers.
9. The sum of the three prominent peaks obtained from each HNGD spectra is called spectral energy profile. The spectral energy profile is convolved with a Gaussian filter of size, average pitch period of that segment. A Gaussian filter of length L is given by

$$G[n] = \frac{1}{\sqrt{2\pi\sigma}} e^{-\frac{n^2}{2\sigma^2}}, n = 1, 2, \dots, L \quad (8)$$

The standard deviation σ used in the above formula is $\frac{1}{4^n}$ of the Gaussian filter length.

10. The spurious peaks are eliminated by using following sub steps:

- (a) First, the spurious peaks are eliminated on the basis that the difference between successive peaks should not be less than 2 ms. This is because 2ms is the minimum range of the pitch period. If two successive peaks having a difference of less than 2ms are found, the peak location with less amplitude is removed.
- (b) Two successive peaks bound a negative region between them. This criteria also eliminates some spurious peak locations.

11. The positive peaks in epoch evidence plot represent epoch locations.

Epoch detection using ZTW method is shown in Fig. 2. The angry emotional speech segment is shown in Fig. 2(a) and its differentiated EGG signal is shown in Fig. 2(b). The spectral energy profile obtained from HNGD spectrum of the speech signal using ZTW analysis is plotted in Fig. 2(c). The epoch evidence plot after convolving spectral energy profile with a Gaussian window of 2 m sec is shown in Fig. 2(d). Epoch locations are shown in Fig. 2(e).

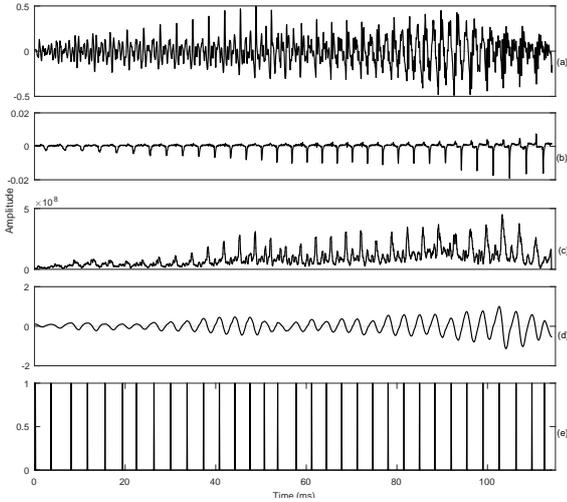


Fig. 2: Epoch extraction using proposed method. (a) Angry speech segment. (b) Differentiated EGG signal. (c) Spectral energy profile obtained from HNGD spectrum. (d) Epoch evidence plot. (e) Epoch locations.

ZTW method for epoch detection is robust for emotional speech [14]. This method is based on spectral peak energy, therefore, it preserves the energy of the signal.

3.3. Epoch Features

The epoch features such as instantaneous pitch, strength of the epoch, slope of the strength of the epoch, the change of phase at the epoch are specific to each emotion [13]. The above mentioned features are determined by the epoch signal obtained by ZTW method [14]. The advantage of this method is that

the value at epoch location is actually the sum of the glottal formants. Therefore, the epochs retain both time and spectral information.

3.3.1. Instantaneous Frequency

Instantaneous Period (IP) is the duration between two successive epoch locations; instantaneous frequency, denoted Δf , is computed as the reciprocal of IP [20, 21]:

$$\Delta f = \frac{1}{t(i) - t(i+1)}, i = 1, 2, \dots, (n-1) \quad (9)$$

where $t(i)$ represents i^{th} epoch location.

3.3.2. Strength Of Excitation

The Strength Of Excitation (SOE) is computed as the difference between two successive epoch values [22]:

$$y(i) = x(i) - x(i+1), i = 1, 2, \dots, (n-1) \quad (10)$$

where $x(i)$ is the epoch strength at i^{th} epoch.

3.3.3. Instantaneous Phase

The instantaneous phase of a glottal signal is obtained by the cosine of the phase function of the corresponding analytical signal.

- The analytic signal $g_a(n)$ corresponding to glottal signal $g(n)$ is given by

$$g_a(n) = g(n) + jg_h(n) \quad (11)$$

- where $g_h(n)$ is the Hilbert transformation of $g(n)$, and is obtained by

$$g_h[n] = IDFT g_h(w) \quad (12)$$

where $g_h(w)$ is defined as:

$$g_h(\omega) = \begin{cases} -jG(\omega), & 0 < \omega < \pi \\ jG(\omega), & -\pi < \omega < 0 \end{cases} \quad (13)$$

$G(\omega)$ is the DTFT of the sequence $g(n)$ and IDFT denotes Inverse Discrete Fourier Transform and

- The Hilbert envelope of glottal signal $g(n)$ is calculated as:

$$h_e[n] = \sqrt{g^2[n] + g_h^2[n]} \quad (14)$$

- The cosine of the phase of the analytic signal $g_a(n)$ is given by

$$\cos\Phi(n) = \frac{Re g_a(n)}{|g_a(n)|} = \frac{g(n)}{h_e[n]} \quad (15)$$

where $g(n)$ is glottal signal derived from speech signal $s(n)$ using ZTW method.

In Fig. 3., instantaneous frequency and SOE values of same speech utterance by same speaker in different emotions are plotted. Figure 3(a) shows instantaneous pitch for two emotions: angry and sad. Red color indicates angry emotion while black indicates Sad emotion. It is clear from Fig. 3(a) that the range of instantaneous pitch varies from 250-400 Hz for angry emotion while for sad it varies from 100-200 Hz. The instantaneous pitch contour for same arousal emotion (happy and angry) is same but their variation with time is different. This property of instantaneous pitch contour is well captured with dynamic model like Hidden Markov Model (HMM) or Long Short Term Memory (LSTM) network. Figure 3(b) shows SOE for two emotions: anger and sad. The variation of SOE is higher in angry emotion than sad emotion. The variation of SOE is quite less in the case of sad emotion. 3(b) shows the phase of glottal signal, it is high for sad compared than angry. The two features SOE and glottal phase also discriminate between same arousal emotion (happy and angry).

4. Development of Emotion Recognition System

Emotion recognition system is an outcome of two principal stages. In the first stage, training is performed using the features extracted from the known emotional speech utterances. In the second stage, i.e., the testing phase, evaluation of the trained model is carried out on unseen emotional speech utterances. The schematic diagram of the proposed emotion recognition system is shown in Fig. 4. We combined the MFCC features with the epoch features namely instantaneous pitch, instantaneous phase and strength of epoch (SOE). The excitation source and system features have complementary information for recognizing emotions, hence, the combined features significantly improve the accuracy of emotion recognition.

4.1. MFCC Feature extraction

Mel Frequency Cepstral Coefficients (MFCCs) features also have emotion specific information. We combine MFCC features with epoch features in our model for recognizing emotions. Gradual spectral variations are captured using 13 MFCCs extracted from speech signal. The speech signal is segmented into frames of size 20 ms, where each frame is overlapped by 10 ms with the adjacent frame. For each frame, 13 MFCC features are extracted. To minimize spectral distortion at the beginning and at the end of each frame, Hamming window is superimposed on each frame segment. MFCC features are extracted from these frames using the MFCC algorithm given in [23]. Recording variations are countered by subtracting cepstral mean and normalizing variance of MFCCs at the utterance level. The schematic diagram of the proposed feature extraction and transformation is shown in Fig. 5.

4.2. DNN-HMMs

In our work, the emotion recognition system has been developed using Hidden Markov model (HMM) [24] – a dynamic modeling approach. It captures the temporal dynamic characteristics of different epoch features of corresponding emotions.

In conventional HMM, the observation probabilities of HMM states are estimated by Gaussian mixture models (GMMs). The GMMs used in such a conventional HMM are statistically inefficient to model non-linear data in the feature space. Therefore, we have replaced the GMMs with DNN to estimate the observation probabilities of observing input sequence at each state in the training phase. In this work, we have developed four HMMs for four discrete emotions. Emotion label is assigned for an unknown speech utterance using Viterbi algorithm. The procedure for training and recognition of DNN-HMM is followed as mentioned in [[25], [26]]. To the best of our knowledge, this is the first time that such a model is being used in an emotion recognition system.

For providing class labels to DNN, we used a GMM-HMM model with five states for each emotion class. Specifically, for each speech utterance in the training set, viterbi algorithm is applied to find an optimal state sequence. The optimal state sequence is stored in the state-label mapping table, which is used to assign a label to each state. The training speech utterances, combined with their labeled state sequences, are then fed as input to the DNN. The output of the DNN is the posterior probabilities of the 20 output units. The observation probability of each state, denoted $p(i_t|q_t)$, is calculated using Bayes theorem as follows:

$$p(i_t|q_t) = \frac{p(q_t|i_t) * p(i_t)}{p(q_t)} \quad (16)$$

Where $I = (i_1, i_2, \dots, i_T)$ is the input sequence and $p(q_t|i_t)$ is the posterior probability obtained as output from the DNN. During decoding, for an unseen speech utterance, the probability of each emotion is estimated and the utterance is assigned the class whose estimated probability is maximum. $p(q_t)$ is computed from the initial state level alignment of the training set. $p(i_t)$ remains constant because input feature vectors are assumed to be mutually independent.

5. Experimental Results and Discussion

Three models were developed for emotion recognition: using system (MFCCs) features, using source (epoch) features, and by combining MFCC and epoch features. The model on combined features has significantly higher accuracy compared to individual models. The experiments were performed on IEMOCAP and IITKGP:SEHSC databases. However, we have conducted experiments for only four emotions, namely angry, happy, sad and neutral. Three-fourth part of the database is used for training purpose and the rest one-fourth of the database is used for evaluating the model. We have used MATLAB tool for feature extraction and KALDI toolkit [27] for developing the system. For the emotion recognition system developed using MFCC features, 13 MFCCs are extracted from each frame. Cepstral mean variance normalization (CMVN) [28] is performed at utterance level to mitigate the recording variations. We have also taken the derivative and double derivative of the normalized MFCCs as features. Therefore, the total number of MFCC features for each frame is 39. To preserve the contextual information, we have used the triphone model approach

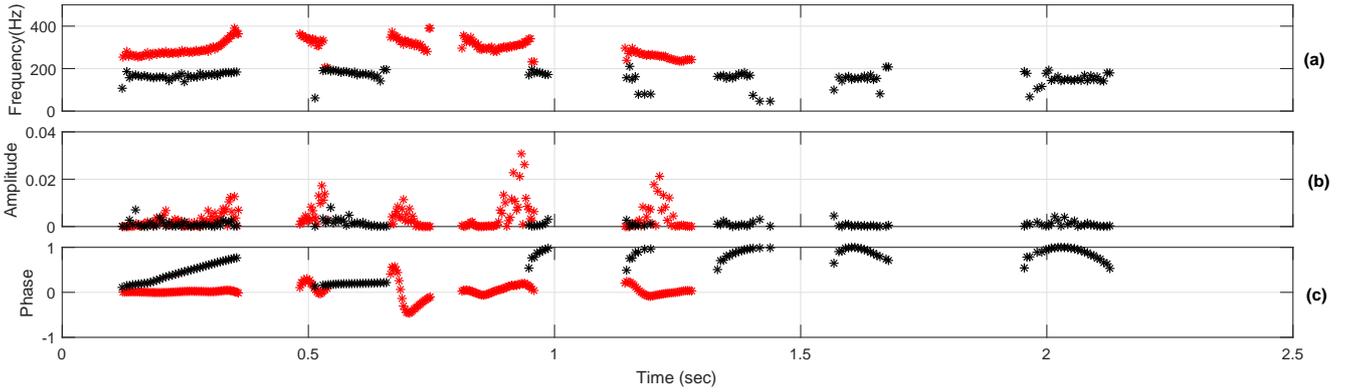


Fig. 3: Instantaneous pitch and SOE contours of angry and sad speech signal using proposed method. (a) Instantaneous pitch contour, (b) SOE contour, and (c) Instantaneous phase contour of angry and sad speech signal.

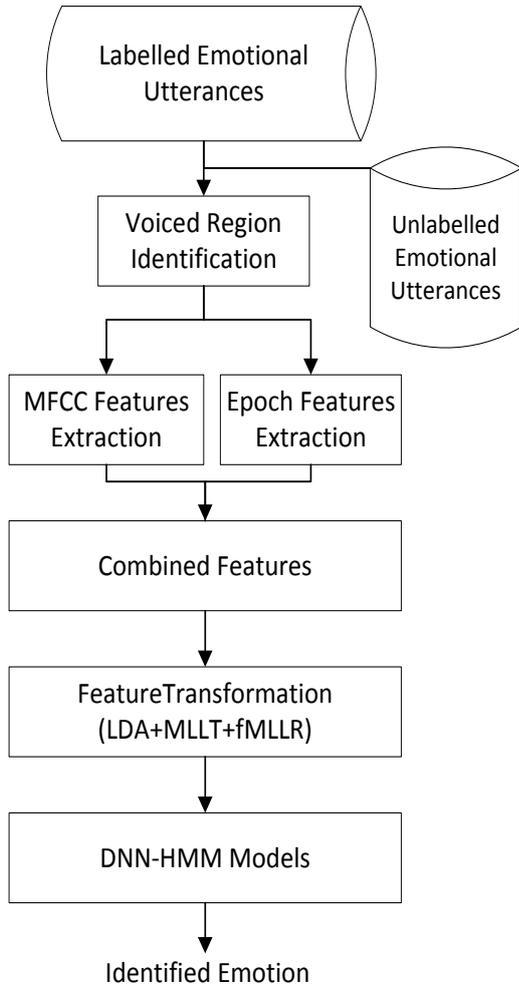


Fig. 4: Schematic diagram of the proposed emotion recognition model.

used in speech recognition where each frame is spliced with the left four frames and the right four frames. A significant improvement in emotion recognition accuracy is observed using the triphone model. Feature transformation is applied on the top of 9 spliced frame features. These features are projected into lower dimensional space using Linear Discriminant Analysis (LDA). Then, diagonalizing Maximum Likelihood Linear Transform (MLLT) [[29, 30]] is applied to further improve the result. Speaker Adaptive Training (SAT) is also used to further enhance the accuracy of the emotion recognition model. For speaker adaptive training Feature Space Maximum Likelihood Linear Regression (fMLLR) transformation is used during both training and testing phases. Thus, accuracy of system is further improved using (LDA+MLLT+SAT)[31]. Four different DNN-HMM models corresponding to each emotion class are built using the transformed feature vectors.

The DNN architecture used is: 80:512x5:20, where 80 is the number of transformed input features to the DNN and 512x5 represents 512 nodes in each of the 5 hidden layers. This DNN configuration was found to be optimal after experimenting with different sized configurations. The results discussed in this paper have been obtained on optimal DNN configuration only. There are 20 output classes in the DNN model (20=4x5, where 4 denotes the number of emotion classes and 5 denotes the number of states in HMM). These output classes are treated as "ground-truth" states and are obtained by GMM-HMM based viterbi algorithm. The initial learning rate of 0.005 is gradually decreased to 0.0005 after 25 epochs. Additional 20 epochs are performed after this. The batch size for training is 512. The training of DNN is performed in three stages as in [32]: (i) unsupervised pre-training consisting in layer-wise training of Restricted Boltzmann Machines (RBM) by Contrastive Divergence algorithm; (ii) frame classification training based on mini-batch Stochastic Gradient Descent (SGD), optimizing frame cross-entropy; and (iii) sequence discriminative training consisting in SGD with per-sentence updates, optimizing state Minimum Bayes Risk (MBR).

In our study, we have considered four categorical (class) la-

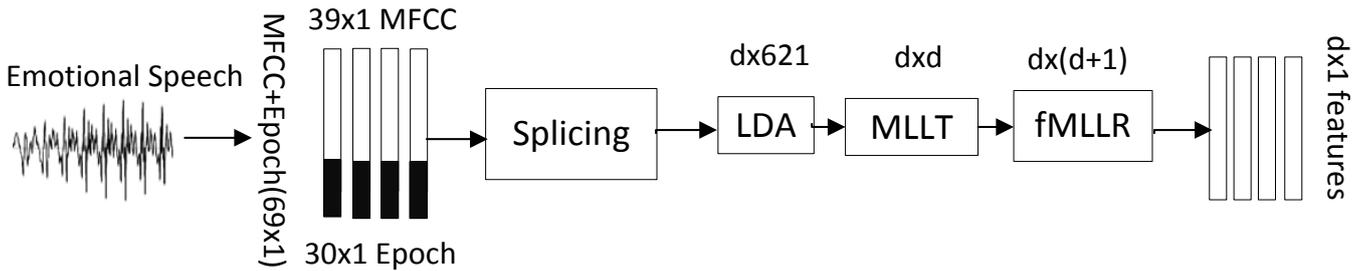


Fig. 5: Schematic diagram of the proposed feature extraction and transformation.

beled emotions namely angry, happy, sad and neutral. The numbers of utterances in each class are 1103, 595, 1084 and 1708 respectively with a total of 4490. The IEMOCAP database is imbalanced. The model was trained in a speaker independent fashion. We used four sessions as training data and the remaining one session for testing. We followed the approach of leave-one-speaker-out cross-validation to generalize the model. The test dataset is also imbalanced corresponding to the emotion classes, hence, we calculated both weighted accuracy(WA) and unweighted accuracy(UWA). Weighted accuracy is calculated by dividing the total number of correct classified test examples with the total number of test samples. Unweighted accuracy is calculated for each emotion category and the average accuracy of all emotions class is taken. The unweighted accuracy is also called class accuracy.

Similarly, for epoch features, the emotion recognition system is developed using three epoch features namely instantaneous pitch, phase and the strength of epoch. These features are extracted using ZTW method. We have taken frames of size 20 ms – same as MFCC features – to extract epoch features. The number of epoch features are different for each frame. To fix the length of epoch-feature vector, we have taken length as 10 – the maximum number of epochs encountered in any frame. If the size of the feature vector is less than 10, we pad the remaining length with zeros. There are no adverse effects of padding to train the network because we transform the input feature vectors (using LDA+MLLT). Therefore, the total number of epoch feature per frame is 30 (10 epochs \times 3 features per epoch). We developed the DNN-HMM model for each emotion using these 30 epoch features.

Finally, we combined the epoch and MFCC features to improve the performance of emotion recognition system. After combining the MFCC and epoch features, the length of the feature vector becomes 69.

We have developed baseline GMM-HMM system using (1) monophone training, (2) triphone training with $MFCC + \Delta + \Delta^2$, and (3) triphone training with LDA+MLLT. We developed the DNN-HMM system with LDA+MLLT. In Table 1 we have shown the result of emotion recognition system using only MFCC and its derivative features. We have also applied LDA+MLLT transformation on MFCC and its derivative features. Our system is trained using both monophone

and triphone training. Triphone system gives better result than monophone because it captures the contextual information. We also estimate the observation probability using DNN instead of GMM as described in previous section. Our system gives best results in the case of DNN-HMM. The average accuracy increases approximately 3.5% when observation probability of HMM models is calculated by DNN instead of GMM. The confusion matrix for experiments done using only $MFCC + \Delta + \Delta^2$ features with LDA+MLLT transformation on DNN-HMM system is shown in Table 2. From the result it is clear that there is more confusion between angry and happy emotions because both are high arousal emotions. The sad and neutral emotions also show confusion because both are low arousal emotions.

Table 1: Emotion classification performance (%) using the MFCC features on IEMOCAP database

| Features | Model | UWA (%) |
|---------------------------------------|---------|----------|
| MFCC(monophone) | GMM-HMM | 44.70 |
| $MFCC + \Delta + \Delta^2$ (triphone) | GMM-HMM | 47.70 |
| MFCC(LDA+MLLT) | GMM-HMM | 51.25 |
| MFCC(LDA+MLLT) | DNN-HMM | 54.35 |

Table 2: Emotion recognition performance on IEMOCAP Database, based on MFCC feature vector of voiced region using DNN-HMM. Abbreviations: A-Anger, H-Happy, N-Neutral, S-Sad

| | MFCC feature vector(Average: 59.58) | | | |
|---------|-------------------------------------|--------------|--------------|--------------|
| | A | H | N | S |
| Anger | 60.21 | 23.29 | 9.45 | 7.05 |
| Happy | 26.56 | 58.17 | 8.70 | 7.57 |
| Neutral | 8.13 | 11.43 | 59.71 | 20.73 |
| Sadness | 8.3 | 8.45 | 23.00 | 60.25 |

Similarly, we also developed the system for epoch features. The average recognition rate for the model developed using MFCC features only is 54.35%. The average recognition rate for the model developed using epoch features only is 54.15 %.

The confusion matrix in Table 3 shows the recognition performance for each emotions using Epoch features. The diagonal elements of the confusion matrix shows the recognition performance for individual emotions using epoch features. From

Table 3: Emotion recognition performance on IEMOCAP Database, based on Epoch feature vector of voiced region. Abbreviations: A-Anger, H-Happy, N-Neutral, S-Sad

| | Epoch feature vector(Average: 54.52) | | | |
|---------|--------------------------------------|--------------|--------------|--------------|
| | A | H | N | S |
| Anger | 57.21 | 15.29 | 22.45 | 5.05 |
| Happy | 13.56 | 52.24 | 21.70 | 12.5 |
| Neutral | 15.23 | 14.40 | 53.71 | 16.66 |
| Sadness | 7.00 | 9.05 | 29.00 | 54.95 |

experimental result it is clear that epoch features discriminate well between angry and happy emotions compared to MFCC features. The average recognition rate for the model developed using the combination of MFCC and epoch features is 60.14%. The performance of the model for each emotion using MFCC features, epoch features and combination of MFCC and epoch features is compared in Table 4. The combined features significantly improves the accuracy of emotion recognition.

Table 4: Emotion classification performance (%) using the Epoch, MFCC and Combined(MFCC+Epoch) features on IEMOCAP database

| Features | Model | UWA (%) |
|-----------------------------------------------------|---------|--------------|
| Epoch Features+LDA+MLLT(triphone) | GMM-HMM | 50.25 |
| | DNN-HMM | 54.15 |
| Epoch Features+LDA+MLLT(triphone) | GMM-HMM | 57.25 |
| Epoch Features+MFCC+ $\Delta + \Delta^2$ (LDA+MLLT) | DNN-HMM | 60.14 |
| Epoch Features+MFCC+ $\Delta + \Delta^2$ (LDA+MLLT) | | |

5.1. Speaker Adaptation

Adaptation is a necessary task for emotion recognition. In general, we train our model with limited dataset but in real environment there may be different types of speakers and noise. There must be a robust method to adapt trained model in real environment. In this paper, we have applied Cepstral mean variance normalization (CMVN) at utterance level to mitigate the recording variations. fMLLR transformation is applied per speaker to adapt the emotion variation of different speakers. After the LDA-MLLT transformation of a feature vector, we transform this matrix into feature space using constraint maximum likelihood linear regression(CMLLR). The model is developed using the strategy of leave-one-speaker-out cross-validation where each time two speakers – that were not a part of the training dataset – are used for testing. There is significant improvement in recognition rate after applying speaker adaptive training for MFCC features. It is mentioned in the Table 5 that after applying fMLLR the emotion recognition rate increases up to 4 % for MFCC features but there is no improvement for epoch features. Therefore, we can say that

Table 5: Emotion classification performance (%) using the Epoch, MFCC and Combined(MFCC+Epoch) features on IEMOCAP database

| Features | Model | UWA(%) |
|--------------------------|---------|--------------|
| MFCC(LDA+MLLT) | DNN-HMM | 54.35 |
| Epoch(LDA+MLLT) | DNN-HMM | 54.15 |
| MFCC+Epoch(LDA+MLLT) | DNN-HMM | 60.14 |
| MFCC(LDA+MLLT+SAT) | DNN-HMM | 59.58 |
| Epoch (LDA+MLLT+SAT) | DNN-HMM | 54.52 |
| MFCC+Epoch(LDA+MLLT+SAT) | DNN-HMM | 64.20 |

epoch features are speaker independent features for which no speaker adaptive technique is required. The bar graph in Fig. 6 Shows that emotion recognition accuracy is higher for combined (MFCC+Epoch) set of features than using each feature-set alone. The average performance of combined features is increased by 5.34% compared to the emotion recognition model developed using MFCC features only. This result proves that both system features and excitation source features have complementary information for emotion recognition.

We evaluate our proposed approach on two databases IEMOCAP and IITKGP:SEHSC database. The bar graph corresponding to MFCC, Epoch and combined feature set for each database is shown in Fig.7. In both databases, the accuracy increased for combined features. The accuracy is more in IITKGP:SEHSC database because it is a scripted database. IEMOCAP database contains both scripted and spontaneous sessions, and is more natural. It is also text independent database. The utterance length in IITKGP:SEHSC database is almost equal whereas large variation in utterance length is observed in IEMOCAP database.

We compare our result with the prior work on IEMOCAP database. In [33], DNN was used to extract the features from speech segment and further utterance level features were constructed and fed to the Extreme Learning Machine(ELM). In [34], raw spectrogram and Low Level Descriptors(LLDs) features were modeled with attentive LSTM. The accuracy is more for LLDs compared to spectrogram. In [35] Convolutional Neural network(CNN) was used for feature extraction from speech frame and these features was fed to the dense neural network. [36], Long Short Term Memory (LSTM) network was used to preserve the contextual information of CNN-based features. CNN-based features are fully data driven. It extracts features from the raw spectrogram which are the representation of speech but it does not contain temporal resolution properly. To achieve temporal resolution, we have to restrict frequency resolution. All the methods used spectrogram for speech representation but it can mislead the accuracy. Our feature extraction approach is not data driven, we are identifying desired temporal and spectral features using signal processing technique. The HMM model captures the contextual information of epoch features. As can be seen from the Table 6 both the weighted and unweighted accuracy outperform from the other methods. Our result proves that MFCC and source features(epoch) contain complimentary information.

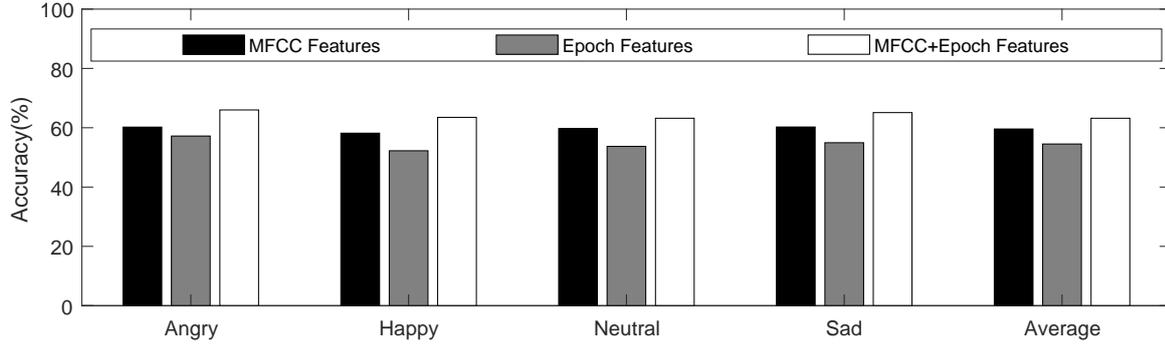


Fig. 6: Emotion classification performance (%) using the Epoch, MFCC and Combined(MFCC+Epoch) features on IEMOCAP database

Table 6: Comparison of Emotion classification performance (%) reported in the prior work on IEMOCAP database

| Model | Features | WA(%) | UWA (%) |
|--------------------------|-----------------------------------------------------------|-------|---------|
| DNN+ELM [33] | MFCC features, pitch-based features and their derivatives | 54.3 | 48.00 |
| LSTM with attention [34] | Local level descriptor and Spectrogram | 63.5 | 58.8 |
| CNN [35] | Spectrogram | 64.78 | 60.89 |
| CNN+LSTM [36] | frame-level Spectrogram | 68.8 | 59.4 |
| DNN-HMM | Epoch (Proposed) | 58.60 | 54.52 |
| DNN-HMM | MFCC (proposed) | 64.3 | 59.58 |
| DNN-HMM | MFCC+Epoch (Proposed) | 69.5 | 64.2 |

6. SUMMARY AND CONCLUSION

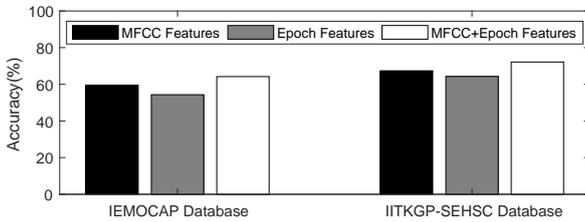


Fig. 7: Emotion classification performance (%) using the Epoch, MFCC and Combined(MFCC+Epoch) features on IEMOCAP and IITKGP-SEHSC databases

The paper highlights the robust characteristic of ZTW method for extracting epoch features. The DNN-HMM model is developed for each emotion using epoch features such as instantaneous pitch, strength of epoch (SOE). The average emotion recognition rate of the proposed model using epoch features is 54.52%. The model developed using epoch features is further combined with the model developed using MFCC feature vectors. The observed accuracy of the proposed model using MFCC and epoch features together is 64.20%. The experimental results show that the epoch feature set is complementary to the MFCC feature set for emotion classification. Our future work is to use LSTM network to capture the contextual information of epoch feature and to explore epoch features in the other applications of speech processing such as speaker identification, speech recognition and, synthesis and language identification.

References

References

- [1] R. Nakatsu, J. Nicholson, N. Tosa, Emotion recognition and its application to computer agents with spontaneous interactive capabilities, *Knowledge-Based Systems* 13 (7) (2000) 497–504.
- [2] D. Ververidis, C. Kotropoulos, A state of the art review on emotional speech databases, in: *Proceedings of 1st Richmedia Conference*, Citeseer, 2003, pp. 109–119.
- [3] T. Sagar, Characterisation and synthesis of emotions in speech using prosodic features, Ph.D. thesis, Masters thesis, Dept. of Electronics and communications Engineering, Indian Institute of Technology Guwahati (2007).
- [4] C. M. Lee, S. S. Narayanan, Toward detecting emotions in spoken dialogs, *Speech and Audio Processing*, *IEEE Transactions on* 13 (2) (2005) 293–303.
- [5] K. E. B. Ooi, L.-S. A. Low, M. Lech, N. Allen, Early prediction of major depression in adolescents using glottal wave characteristics and teager energy parameters, in: *Acoustics, Speech and Signal Processing (ICASSP)*, 2012 IEEE International Conference on, IEEE, 2012, pp. 4613–4616.
- [6] L.-S. A. Low, N. C. Maddage, M. Lech, L. B. Sheeber, N. B. Allen, Detection of clinical depression in adolescents speech during family interactions, *IEEE Transactions on Biomedical Engineering* 58 (3) (2011) 574–586.
- [7] Y. Yang, C. Fairbairn, J. F. Cohn, Detecting depression severity from vocal prosody, *IEEE Transactions on Affective Computing* 4 (2) (2013) 142–150.
- [8] C. Clavel, I. Vasilescu, L. Devillers, G. Richard, T. Ehrette, Fear-type emotion recognition for future audio-based surveillance systems, *Speech Communication* 50 (6) (2008) 487–503.
- [9] Y. Wang, L. Guan, An investigation of speech-based human emotion recognition, in: *Multimedia Signal Processing*, 2004 IEEE 6th Workshop on, IEEE, 2004, pp. 15–18.
- [10] J. Nicholson, K. Takahashi, R. Nakatsu, Emotion recognition in speech using neural networks, *Neural computing & applications* 9 (4) (2000) 290–296.
- [11] D. Ververidis, C. Kotropoulos, I. Pitas, Automatic emotional speech classification, in: *Acoustics, Speech, and Signal Processing*, 2004. *Proceedings. (ICASSP'04)*. IEEE International Conference on, Vol. 1, IEEE, 2004, pp. 1–593.
- [12] F. Dellaert, T. Polzin, A. Waibel, Recognizing emotion in speech, in: *Spoken Language*, 1996. *ICSLP 96. Proceedings.*, Fourth International Conference on, Vol. 3, IEEE, 1996, pp. 1970–1973.
- [13] S. R. Krothapalli, S. G. Koolagudi, Characterization and recognition of emotions from speech using excitation source information, *International journal of speech technology* 16 (2) (2013) 181–201.
- [14] J. Yadav, M. S. Fahad, K. S. Rao, Epoch detection from emotional speech signal using zero time windowing, *Speech Communication*.
- [15] C. Busso, M. Bulut, C.-C. Lee, A. Kazemzadeh, E. Mower, S. Kim, J. N. Chang, S. Lee, S. S. Narayanan, Iemocap: Interactive emotional dyadic motion capture database, *Language resources and evaluation* 42 (4) (2008) 335.
- [16] S. G. Koolagudi, R. Reddy, J. Yadav, K. S. Rao, Iitkgp-sehsc: Hindi speech corpus for emotion analysis, in: *Devices and Communications (ICDeCom)*, 2011 International Conference on, IEEE, 2011, pp. 1–5.
- [17] S. S. Kumar, K. S. Rao, Voice/non-voice detection using phase of zero frequency filtered speech signal, *Speech Communication* 81 (2016) 90–103.
- [18] Y. Bayya, D. N. Gowda, Spectro-temporal analysis of speech signals using zero-time windowing and group delay function, *Speech Communication* 55 (6) (2013) 782–795.
- [19] K. S. R. Murty, B. Yegnanarayana, Epoch extraction from speech signals, *IEEE Transactions on Audio, Speech, and Language Processing* 16 (8) (2008) 1602–1613.
- [20] S. G. Koolagudi, R. Reddy, K. S. Rao, Emotion recognition from speech signal using epoch parameters, in: *Signal Processing and Communications (SPCOM)*, 2010 International Conference on, IEEE, 2010, pp. 1–5.
- [21] N. Narendra, K. S. Rao, Robust voicing detection and f_0 estimation for hmm-based speech synthesis, *Circuits, Systems, and Signal Processing* 34 (8) (2015) 2597–2619.
- [22] P. Gangamohan, S. R. Kadiri, S. V. Gangashetty, B. Yegnanarayana, Excitation source features for discrimination of anger and happy emotions., in: *INTERSPEECH*, 2014, pp. 1253–1257.
- [23] L. Rabiner, B.-H. Juang, *Fundamentals of speech recognition*.
- [24] L. Rabiner, B. Juang, An introduction to hidden markov models, *ieee assp magazine* 3 (1) (1986) 4–16.
- [25] L. Li, Y. Zhao, D. Jiang, Y. Zhang, F. Wang, I. Gonzalez, E. Valentin, H. Sahli, Hybrid deep neural network–hidden markov model (dnn-hmm) based speech emotion recognition, in: *Affective Computing and Intelligent Interaction (ACII)*, 2013 Humaine Association Conference on, IEEE, 2013, pp. 312–317.
- [26] G. Hinton, L. Deng, D. Yu, G. E. Dahl, A.-r. Mohamed, N. Jaitly, A. Senior, V. Vanhoucke, P. Nguyen, T. N. Sainath, et al., Deep neural networks for acoustic modeling in speech recognition: The shared views of four research groups, *IEEE Signal Processing Magazine* 29 (6) (2012) 82–97.
- [27] D. Povey, A. Ghoshal, G. Boulianne, L. Burget, O. Glembek, N. Goel, M. Hannemann, P. Motlicek, Y. Qian, P. Schwarz, et al., The kaldi speech recognition toolkit, in: *IEEE 2011 workshop on automatic speech recognition and understanding*, no. EPFL-CONF-192584, IEEE Signal Processing Society, 2011.
- [28] O. Viikki, K. Laurila, Cepstral domain segmental feature vector normalization for noise robust speech recognition, *Speech Communication* 25 (1-3) (1998) 133–147.
- [29] M. J. Gales, Maximum likelihood linear transformations for hmm-based speech recognition, *Computer speech & language* 12 (2) (1998) 75–98.
- [30] M. J. Gales, Semi-tied covariance matrices for hidden markov models, *IEEE transactions on speech and audio processing* 7 (3) (1999) 272–281.
- [31] S. P. Rath, D. Povey, K. Vesely, J. Cernocký, Improved feature processing for deep neural networks., in: *Interspeech*, 2013, pp. 109–113.
- [32] K. Vesely, A. Ghoshal, L. Burget, D. Povey, Sequence-discriminative training of deep neural networks., in: *Interspeech*, 2013, pp. 2345–2349.
- [33] K. Han, D. Yu, I. Tashev, Speech emotion recognition using deep neural network and extreme learning machine, in: *Fifteenth Annual Conference of the International Speech Communication Association*, 2014.
- [34] S. Mirsamadi, E. Barsoum, C. Zhang, Automatic speech emotion recognition using recurrent neural networks with local attention, in: *Acoustics, Speech and Signal Processing (ICASSP)*, 2017 IEEE International Conference on, IEEE, 2017, pp. 2227–2231.
- [35] H. M. Fayek, M. Lech, L. Cavedon, Evaluating deep learning architectures for speech emotion recognition, *Neural Networks* 92 (2017) 60–68.
- [36] A. Satt, S. Rozenberg, R. Hoory, Efficient emotion recognition from speech using deep learning on spectrograms, *Proc. Interspeech 2017* (2017) 1089–1093.