CrossMark

# An Automatic Tamil Speech Recognition system by using Bidirectional Recurrent Neural Network with Self-Organizing Map

S. Lokesh[1] · Priyan Malarvizhi Kumar[2] · M. Ramya Devi[3] · P. Parthasarathy[2] · C. Gokulnath[2]

## Abstract

Speech recognition is one of the entrancing fields in the zone of computer science. Exactness of speech recognition framework may decrease because of the nearness of noise exhibited by the speech signal. Consequently, noise removal is a fundamental advance in automatic speech recognition (ASR) system. ASR is researched for various languages in light of the fact that every language has its particular highlights. Particularly, the requirement for ASR framework in Tamil language has been expanded broadly over the most recent couple of years. In this work, bidirectional recurrent neural network (BRNN) with self-organizing map (SOM)-based classification scheme is suggested for Tamil speech recognition. At first, the input speech signal is pre-prepared by utilizing Savitzky–Golay filter keeping in mind the end goal to evacuate the background noise and to improve the signal. At that point, Multivariate Autoregressive based highlights by presenting discrete cosine transformation piece to give a proficient signal investigation. And in addition, perceptual linear predictive coefficients likewise separated to enhance the classification accuracy. The feature vector is shifted in measure, for picking the right length of feature vector SOM utilized. At long last, Tamil digits and words are ordered by utilizing BRNN classifier where the settled length feature vector from SOM is given as input, named as BRNN-SOM. The experimental analysis demonstrates that the suggested conspire accomplished preferable outcomes looked at over exist deep neural network–hidden Markov model algorithm regarding signal-to-noise ratio, classification accuracy, and mean square error.

**Keywords** Automatic Tamil Speech Recognition · Preprocessing · Feature extraction · Classification · Bidirectional Recurrent Neural Network (BRNN) · Self-Organizing Map (SOM) · Savitzky–Golay Filter (SGF) · Multivariate Autoregressive (MAR) · Discrete Cosine Transformation (DCT) · Perceptual Linear Predictive (PLP)

# 1 Introduction

Digital speech is normally the most agreeable method of interaction in the field of human–computer interaction (HCI). Voice recognition is an undertaking of translating human speech into a digitized type of speech that can be deciphered by gadget of a PC. Automatic speech recognition (ASR) framework winds up noticeably difficult

because of different kinds of speaker, talking style, environment, noise, and so on. In spite of its impediments, speech recognition innovation is a significant device in numerous applications like live subtitling on TV, correspondence in medical interpretations, command control in robotics, speech-to-text transformation for note making frameworks, and substitution of keyboard and mouse for physically or outwardly tested individuals. ASR is a procedure by which a machine recognizes discourse. It takes a human expression as an input and provides a series of words as result. Such research on ASR frameworks is basically created for the English language; however, for Indian languages, it is still in prior stage. Tamil language is one of the broadly spoken languages in the world with more than 77 million speakers. Thus, there is a pressing requirement for the framework to communicate with Tamil language.

✉ Priyan Malarvizhi Kumar
  priyanm.k2016@vitstudent.ac.in

1 Department of Computer Science and Engineering, Hindustan Institute of Technology, Coimbatore, India

2 VIT University, Vellore, India

3 Department of Computer Science and Engineering, Hindusthan College of Engineering and Technology, Coimbatore, India

RETRACTED ARTICLE

In noisy environment, the exactness level of ASR framework will endure significantly [1]. Though execution of ASR has seen general enhancements, the comparative debasement within the sight of noise or resonation keeps on being a significant test in the growing real-world applications [2]. The answer for beating the execution debasement in noisy environment is the utilization of multi-condition training data set [3], where the acoustic models were prepared utilizing the information from the objective space. Be that as it may, in a sensible situation, it is not generally conceivable to get sensible measures of training information from a wide range of noisy conditions. In multi-condition training, the execution of ASR frameworks is essentially more awful contrasted with noisy free or clean conditions. The main objective of this paper is to mention the robustness in feature extraction stage of ASR.

As of late, other experiments were examined by utilizing neural networks (NN) to take in the nonlinear mapping among perfect and partial speech feature coefficients. The neural systems are widespread, that could be utilized for the two major problems such as classification and regression. The strategy of NN has been effectively utilized to enhance ASR [4]. A NN is utilized with more than one hidden layer is normally known as deep NN, or DNN. As of late, DNN has turned out to be prominent after a pre-training step, known as restrictive Boltzmann machine (RBM) pre-training [5, 6], was acquainted with instate the system parameters to some sensible esteems to such an extent that back engendering would then be able to be utilized to prepare the system proficiently on task-dependent target capacities. The main advantage of the DNN with more than one hidden layer is that the profound model of the DNN permits significantly increases productive portrayal of numerous nonlinear transformations [7]. The DNN and other neural systems have been connected to numerous speech-processing tasks. The DNN was developed for providing acoustic displaying in ASR frameworks, and now it turned into the accepted typical acoustic model [8]. In [24], DNN design, known as the deep recurrent neural network (DRNN), is utilized to obtain the clean speech features using MFCC from noisy speech features. An uncommon instance of RNNs, known as long short-term memory (LSTM), is utilized to map reverberant features to clean features in [9]. In [10], a DNN is utilized to foresee the speech cover, which is utilized for upgrading speech for robust ASR system. It is likewise discovered that adjusting the mask-estimating DNN utilizing direct information change additionally enhances the ASR execution, while the past two examinations concentrate on anticipating low-dimensional feature extraction for ASR. In [11], DNN is used specifically to gauge the high-measurement log-size range for speech de-noising. The same strategy was used to connect in future in the preprocessor

for ASR system [12]. The analysis utilizes NN with outfit classifier to evaluate low-dimensional speech feature extraction for the ASR undertaking and a high-dimensional log-magnitude spectrum extraction for the speech upgrade assignment.

The fundamental target of this paper is to actualize the classification and recognition system for Tamil spoken words. To recoup unique speech from noisy speech signal, the preprocessing plan is finished by utilizing SGF. Feature extraction (FE) is one of the huge strides in ASR framework which changes original signal into a shape that is fitting for the classification model. To complete this undertaking, two imperative features like MAR and PLP are separated for effectual classification. BRNN is straightforward nonlinear classifier and has bigger adaptability in taking care of classification task. The procedure of feature extraction brings about variable length of feature vector for every spoken word. To change over-factor estimate feature vector into fixed-size feature vector, SOM is utilized as a contribution to be encouraged into the ensemble classifier [13]. The experimental analysis demonstrates that the proposed plot achieved preferred outcomes when distinguished with other plans.

This paper is organized as follows: Sect. 2 explains about few related works in Tamil speech recognition, and Sect. 3 describes various methodologies of the proposed system. Section 4 provides the recognition results of experiments with and without de-noising procedure. At last, conclusion is given in Sect. 5.

## 2 Related work

Here, a portion of the related works in Tamil speech acknowledgment was given. Radha et al. [14] took a shot at separated words for Tamil spoken language, and here input signal was preprocessed utilizing four sorts of filters, and from best filter output, LPCC feature extraction was finished. The classification and recognition received utilizing back-propagation neural system, which has created better outcomes for restricted vocabulary.

Radha et al. [15] exhibited a continuous speech recognition (CSR) framework for Tamil language with the help of hidden Markov model (HMM). In feature extraction, MFCC feature extraction is utilized as a preprocessing stage or front-end for the proposed framework. The monophone-based acoustic model is perceived to give the arrangement of sentences from medium vocabulary. The outcomes are observed to be acceptable with word recognition accuracy of 92 and 81% of sentence exactness for the proposed framework.

Patel and Rao [16] proposed the traditional approach; low recurrence MFCC vectors are removed and grilled

with recurrence sub-band decomposition. The executed framework indicates preferred productivity over-existing MFCC technique. Chandrasekar and Ponnavaikko [17] built up a speaker subordinate consistent speech recognition framework for Tamil. The proposed strategy portions words from sentences and afterward character from words. The back-propagation algorithm is utilized for training and verifying a framework. The framework was tried for sectioning words from nine spoken sentences and accomplishes precision of 80.95%.

Rojathai and Venkatesulu [18] displayed the novel speech word acknowledgment framework for Tamil which comprises of three phases. The primary input speech signal is pre-prepared utilizing Gaussian filtering procedure. From noiseless flag, MFCC feature vectors were extricated from training dataset and test dataset. At that point, feed-forward back-propagation neural network (FFBNN) experiences training and testing with their particular datasets. The execution of proposed method provides preferred acknowledgment result over-existing HMM and associative ANN system.

Sigappi and Palanivel [19] detailed a speaker-dependent medium-sized vocabulary Tamil speech recognition mechanism. Here the framework was prepared and tried with HMM and auto-associative neural networks (AANN) utilizing 8000 and 2000 examples individually. The MFCC feature extraction procedures were connected to input speech tests to extricate feature vectors. The execution expresses that HMM with five states and four blends yields high-acknowledgment execution than AANN.

Sivaraj and Rama [20] proposed the speaker-independent isolated Tamil words recognition framework utilizing discrete wavelet transform (DWT) and multilayer perceptron arrange prepared with back-propagation training algorithm. The db4 sort of wavelet utilized for wavelet-based feature extraction. At that point, the speech tests in database progressively experience an eight-level disintegration to get estimate and detail coefficients. Here 70% of information is utilized for training, 15% for approval, 15% for testing, and at the end, it accomplishes general acknowledgment precision of 90%.

Chandrarajan et al. [21] built up the speaker-independent triphone-based medium vocabulary-persistent speech recognizer for Tamil language. The usage of the framework is finished with Sphinx-4 structure of HMM show with three discharging states and one non-emitting state with nonstop thickness of 8 Gaussian per state was utilized. They built a phoneme-based context-dependent acoustic model for 1700 remarkable words, at that point pronunciation dictionary with 44 base telephones and triphone-based measurable language model. The framework brings about great word precision and same word blunder rate for training and test expressions.

Saraswathi and Geetha [22] enhanced the precision of Tamil speech framework by planning language models at different levels such as segmentation phase, recognition phase, syllable, and word level error correction phase. They enhanced the acknowledgment precision at each stage, and lastly 87.1% exactness was acquired. Karpagavalli et al. [23] created speaker-independent isolated Tamil digits recognition utilized and accomplished general acknowledgment exactness of 91.8%. From input discourse signals, MFCC feature vectors were removed and prepared utilizing vector quantization (VQ) approach. The codebook for every digit is produced utilizing Linde–Buzo–Gray (LBG) VQ training algorithm. Iswarya and Radha [24] outlined the system for Tamil speech-based query processing design to recovery English textual documents. They coordinated speech recognition and cross-language content recovery framework.

From a few related works of Tamil speech recognition, it is discovered that a significant number of the exploration were performed with the help of MFCC, LPC, and wavelet-based feature extraction procedures. Additionally, for recognition purpose, hidden Markov demonstration and neural systems were utilized by many creators. At that point, few papers made utilization of noise-filtering systems for noise evacuation.

# 3 Proposed methodology

In this section, the suggested BRNN-SOM step-by-step process has been explained. Here, the Tamil speech recognition is indicated.
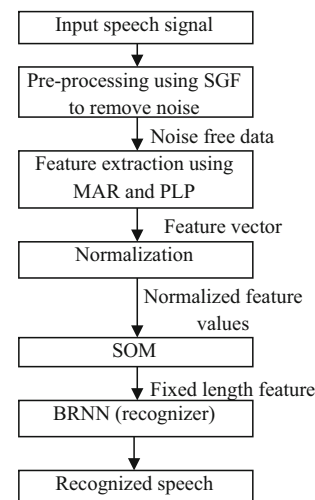


**Fig. 1** Architecture of RNN-SOM based speech recognition system

## 3.1 System overview

The suggested speech recognition system is shown in Fig. 1. The initial phase in speech recognition is pre-preparing speech signals which diminish noise in view of SGF noise removal algorithm. At that point, the MAR- and PLP-based features were removed for effectual classification. To keep away from mutilations in speech signal, cepstral mean standardization method is connected [13]. With a specific end goal to make fixed-length trajectory model input to BRNN classifier, SOM is connected to the feature vectors.

## 3.2 Preprocessing

Pre-preparing of a speech signal is considered as an essential advance in the improvement of a robust speech or a speaker recognition system. To upgrade the precision and productivity of speech recognition framework, speech signals are for the most part pre-prepared before they additionally break down. Here, SGF conspire is utilized to expel white noise from input speech signal. This SG separating chooses the ideal frame size and request utilizing iterative examination and signal relationship. This kills the heuristic view, frequently cited as a drawback, in SG filter. Assist the processing speed is expanded fundamentally. The filter coefficients should be assessed only once for ASR application which influences the filtering process to be basic, simple, and quick. Because of the above reason, we utilized SGF conspire in our framework.

### 3.2.1 Savitzky–Golay filter (SGF)

Ordinarily, this digital filter utilizes the system of linear least squares for data smoothing, which gets a high signal-to-noise ratio and holds the initial state of the signal. With its numerous favourable circumstances over standard filtering systems, SGF is favored for recovering original signal structure while expelling noise in this work.

Savitzky–Golay channel is connected to a series of advanced information focuses on the point of expanding the signal-to-noise ratio without distorting the signal. Obtaining the subsets of successive data points was fitted utilizing a low-request polynomial with linear least-square method, and convolution of the considerable number of polynomials is then acquired [25, 26]. The information having a set of n $\{x_i, y_j\}$ points, where $j = 1, 2 \ldots n$, and x is an independent variable, while y is a observed esteem, can be spoken with an set of m convolution coefficients, $C_i$, and provided as

$$Y_j = \sum_{i=-(m-1)/2}^{i=(m-1)/2} C_i y_{j+i} \quad \frac{m+1}{2} \le j \le n - \frac{m-1}{2} \tag{1}$$

Execution of SG filter typically demands three sources of information: the noisy signal $(x)$, the order of the polynomial $(k)$, and its frame size $(f)$. The best-fit estimations of $k$ and $f$ for a signal are by and large evaluated utilizing experimentation strategy (trial and error method). On the other hand, the qualities can likewise be acquired utilizing prior experience or already assessed values for a specific level of SNR for the provided signal. The filtered signal is acquired and assessed over the range of qualities.

## 3.3 Feature extraction

For the most part, the feature extraction process turns out to be exceptionally troublesome because of different requirements engaged with speech input. They are: (1) speech signal varying for a given word between speakers, (2) replication utterances by a similar speaker, (3) accent difference between speakers. To understand the above impediments, a great feature extraction strategy ought to be equipped for distinguishing particular properties that are more important to the linguistic substance. Additionally, it should dispose of all other insignificant data such as background noise, channel distortion, emotion and so forth. In this manner, the decision of feature extraction turned out to be extremely critical in pattern recognition issue. In this way, to take care of the above issue, here we presented two sorts of feature that were extricated plans namely MAR and PLP coefficient features for useful classification. The MAR strongly worked for noise-free and noisy information, because of the long-haul discrete cosine change.

### 3.3.1 MAR feature extraction

In proposed work, endeavor to mutually demonstrate the transient covers the various subgroups utilizing a time series approach [27, 28]. The multivariate AR (MAR) demonstrating procedure is one of the strategies for approximating the random time series vector as a linear combination of "past" vectors. In this strategy, the forecast coefficients are evaluated by utilizing the generalized least squares. In this, MAR modeling is broadly utilized as a part of econometrics for anticipating applications [29]. This investigation speaks to the primary use of MAR modeling utilizing multi-band Riesz to observe the best estimation.

To improve the application of speech processing, it utilizes the discrete cosine transform (DCT) coefficients of different spectral groups in the MAR system. Generally, MAR modeling protects the peak signals in the joint spectro-temporal domain and endeavors the 2D structure of

speech spectrograms. Provided with the absence of time–frequency connections in noisy environment, this suggested 2D modeling permits the extraction of the multiband features illustrative with basic speech signal even within the sight of noisy condition.

Figure 2 shows the block diagram of the proposed approach for feature extraction. The fragments of the input speech signal vary from 2000 ms of non-overlapping windows, which are changed utilizing DCT. Obtained full-band DCT signal is windowed into a set of 39 overlapping subgroups utilizing Gaussian-shaped windows with center frequencies picked consistently with the mel scale. The obtained windowing is like mel band windowing done in traditional feature extraction such as Mel Frequency Cepstral Coefficients (MFCC). The sequences of DCT numerous sub-bands are loaded together to frame vector series data $y_q$ ($q$ signifies the coefficient index in DCT) is given in Eq. (2).

$$y_q = \sum_{k=1}^{p} A_k y_{q-k} + u_q \tag{2}$$

where $y$ is determined as provided, $D$ is dimensional vector process of sequential data indexed by $q = 1 \ldots Q$, a multivariate AR model of order $p$ is indicated above, and $u$ is a D-dimensional white noise random process with a covariance matrix $\sum u$, and the MAR coefficients $A_k$ are square matrices of size $D$ which characterize in the model.

The procedure of the MAR model estimation is enforced, and model parameters $\beta$ are computed in Eq. (3).

$$\hat{\beta} = \left( (ZZ^T)^{-1} Z \otimes I_k \right) \eta \tag{3}$$

where $\eta = \text{vec}(BZ) + u$, $u = \text{vec}(U)$, $B := [A_1, A_2, \ldots, A_p]$, $U := [u_1, u_2, \ldots, u_Q]$, $Z := [Z_0, \ldots, Z_{Q-1}]$ of dimension $Dp \times Q$, $\otimes$ is the Kronecker product and $I_k$ is the identity matrix of size $k$.

We make use of a fixed model order of $p = 160$ for estimating the MAR of 2000 ms of speech signal. The temporal envelopes of the sub-band are then estimated with the help of Eq. (4).

$$\hat{s}_y[n] = \text{diag}\left( H[n]^{-1} \widehat{\sum u} H[n]^{-1} \right) \tag{4}$$
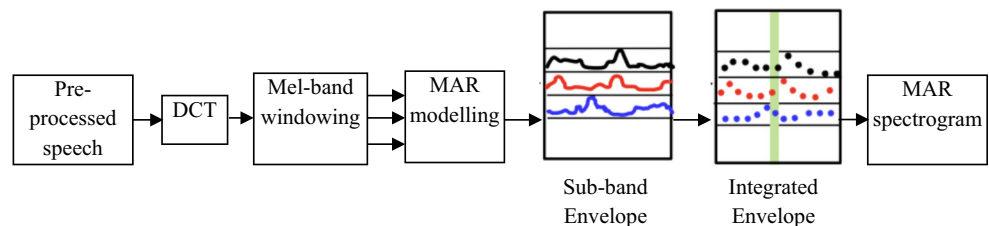
where $S_y[n]$ indicated the Riesz envelope which is an extension of Hilbert envelope to 2D signals of various speech sub-bands. Later the MAR estimate of the Riesz envelope is provided above (for $H[n] = H[z]|_{z=e^{-j2\pi n}}$), where $H[z] = I_D - \sum_{k=1}^{p} A_k z^{-k}$, it is a multidimensional $z$-transform filter. Here, the DCT coefficients of three mel bands are utilized in MAR modeling (i.e., $D = 3$).

The sub-band MAR envelopes are coordinated with a Hamming window over a 25-ms window with a 10-ms move. The combination in time of the sub-band envelopes provides a gauge of the MAR spectrogram of the input signal. The discrimination of the spectrographic portrayal from MAR displaying and the ordinary mel spectrogram is shown in Fig. 3. As observed, the MAR displaying brings about a smooth portrayal, and this underlines just the high vitality locales of the signal. The combined estimation envelopes are required by the 2D spectro-temporal modeling which likewise enables the model to concentrate basically on time–frequency relationships of the fundamental speech signal while suppressing the impacts of noise as delineated by the portrayals obtained for the bauble noise at 10 dB SNR with the presence of channel noise. The properties of the MAR demonstrate enhancing of the noise power in the portrayals obtained from this method. In ASR feature extraction, the incorporated sub-band temporal envelopes for span of 200 ms (centered on a 10-ms outline) changed to 14 coefficients of DCT for every sub-band. The features of MAR are likewise added with spectral delta features yielding 1092 features.

### 3.3.2 Perceptual linear predictive coefficients

Perceptual linear predictive (PLP) shows an optional method to MFCC yet utilized less every now and again. The primary distinction between Mel scale cepstral investigates and PLP is identified with yield cepstral coefficients. PLP alters the transient range of speech more precisely than LPC models, by enforcing few psychophysically based changes. The PLP utilizes an all-pole model to smooth the altered power spectrum, where the yielded cepstral coefficients are then processed in light of this case. In PLP, the spectrum is distorted by the Bark scale filter bank of 18 filters for covering the frequency

**Fig. 2** Block diagram of the MAR spectrogram model



Pre-processed speech → DCT → Mel-band windowing → MAR modelling → Sub-band Envelope → Integrated Envelope → MAR spectrogram
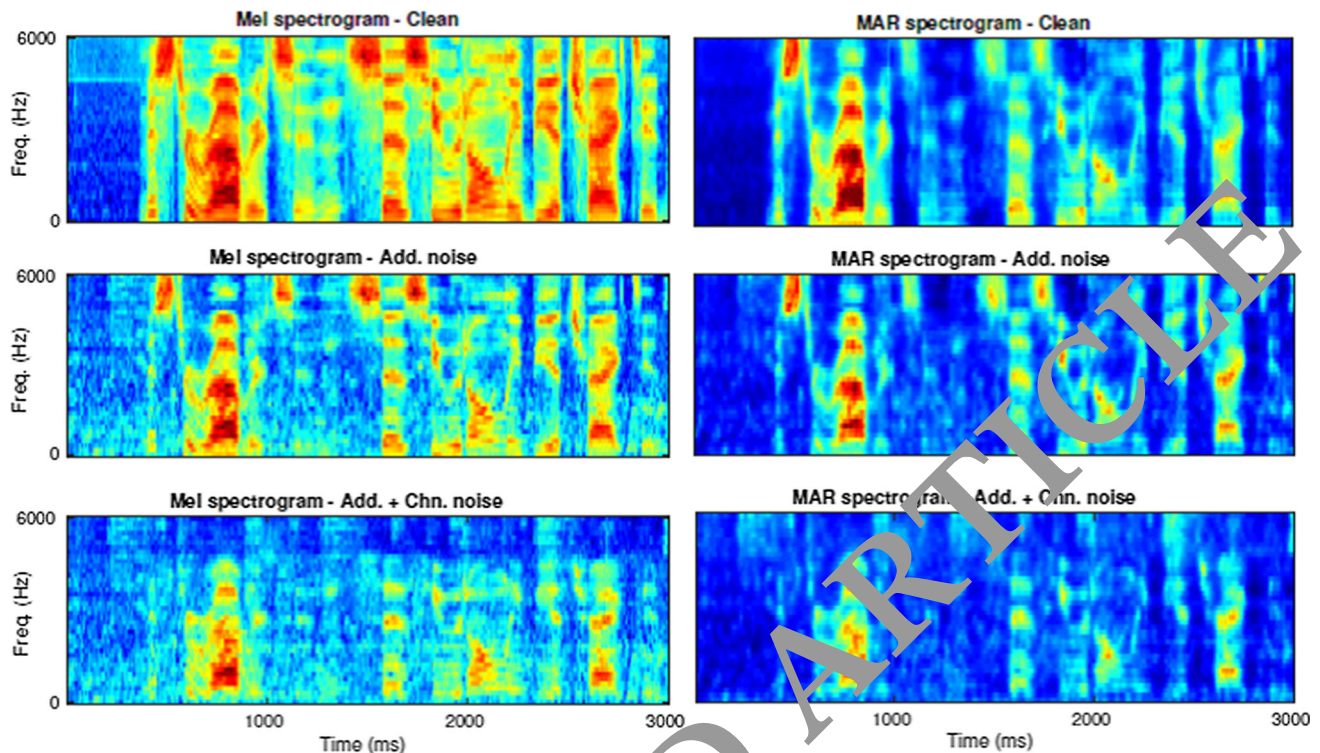
**Fig. 3** Comparison of mel spectrogram estimation using MAR model with conventional mel spectrogram for clean and noisy speech recordings from FIRE dataset

scope of (0, 5000) Hz. The Bark scale is indicated in Eq. (5).

$$\text{Bark}(f) = \frac{26.81f}{1960 + f} - 0.53 \tag{5}$$

The resultant filter bank energies are increased by an equivalent loudness curve [30]. Thus, the critical band filter is figured through discrete convolution of power range with piecewise guess. From that point onward, cube root compression pressure improved the situation yield amplitudes to re-enact the power law of hearing. Then, the equal loudness pre-emphasis is utilized to down-specimen the signal, and an inverse discrete Fourier transform (IDFT) is connected to get equivalent autocorrelation function. At last, PLP coefficients are processed by changing over-autoregressive coefficients to cepstral coefficients. In this examination work, a successful basic band of 24 filter banks is utilized. The MAR and PLP coefficient esteems are standardized by utilizing cepstral mean standardization strategy.

### 3.4 Classification

The procedure feature extraction brings about factor length of feature vectors where SOM is a neural system that proselytes fluctuating size into fixed size of features vectors that will bolster into the classifier as input. At that point,

the utilization of SOM with BRNN enhanced the recognition accuracy and limits the preparation time further. SOM is unsupervised learning technique that works in light of competitive leaning strategy. The SOM algorithm utilizes as input the variable length feature vector and maps it to a steady size of six groups while safeguarding the input size. The algorithm comprises of three undertakings, to be specific: competitive task, cooperative task, and adaptation task. This area presents points of interest of the regular RNN and group classifiers.

In a self-organizing map, the neurons are put at the hubs of a cross section that is generally of one dimension or two dimensions. The higher-dimensional maps are likewise conceivable, however, not as normal. The neurons turn out to be specifically tuned to different input patterns (stimuli) or classes of input patterns over the span of learning process. The areas of neurons so tuned (i.e., the winning neurons) wind up noticeably requested concerning each other such that an important organizing of framework for various input feature is made over the grid.

RNN have feedback associations and address the transient relationship of contributions by keeping up inner states that have memory. RNN are systems with at least one input association. A feedback association is utilized to pass output of a neuron in a specific layer to the past layer(s) [31]. The variation among MLP and RNN will be

RNN have encouraged forward association for all neurons (completely associated). Subsequently, the associations' permits the system demonstrate the dynamic conduct. RNN is by all accounts more normal for speech recognition than MLP in light of the fact that it permits fluctuation in input length [32].

The inspiration for enforcing recurrent neural network to this space is to exploit their capacity to process short-term spectral features yet react to long-term temporal events. Past research has affirmed that speaker acknowledgment execution enhances as the length of expression is expanded [33]. Likewise, it has been demonstrated in ID issues. RNNs may present a superior execution and learn in a shorter time than regular encourage forward systems [34].

Provided an input variable length feature vector sequence $x = (x_1, \ldots, x_T)$, a standard recurrent neural network (RNN) estimates the hidden vector sequence $h = (h_1, \ldots, h_T)$ and output vector sequence $y = (y_1, \ldots, y_T)$ by repeating the following equations from $t = 1$ to $T$:

$$h_t = \mathcal{H}(W_{xh}x_t + W_{hh}h_{t-1} + b_h) \tag{6}$$

$$y_t = W_{hy}h_t + b_y \tag{7}$$

where the $W$ terms indicate the weight matrices (e.g., $W_{xh}$ is the input hidden weight matrix), the $b$ terms indicates bias vectors (e.g., $b_h$ is hidden bias vector) and $H$ is the hidden layer function. $H$ is generally an element-wise application of a sigmoid function.

One inadequacy of customary RNNs is that they are just ready to make utilization of past setting in speech recognition, where entire expressions are deciphered without a moment's delay, there is no reason not to abuse future setting also. Bidirectional RNNs (BRNNs) [35] does this by handling the information in the two headings with two separate hidden layers, which are then encouraged advances to a similar yield layer.

As demonstrated in Fig. 4, a BRNN estimates the forward hidden sequence $\vec{h}_t$, the backward hidden sequence $\overleftarrow{h}_t$, and the output sequence $y$ by repeating the backward layer
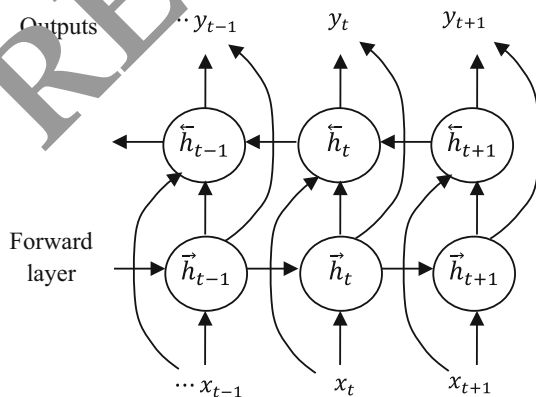


**Fig. 4** Architecture of BRNN

from $t = T$ to 1, the forward layer from $t = 1$ to $T$, and then updating the output layer:

$$\vec{h}_t = \mathcal{H}\left(W_{x\vec{h}}x_t + W_{\vec{h}\vec{h}}\vec{h}_{t-1} + b_{\vec{h}}\right) \tag{8}$$

$$\overleftarrow{h}_t = \mathcal{H}\left(W_{x\overleftarrow{h}}x_t + W_{\overleftarrow{h}\overleftarrow{h}}\overleftarrow{h}_{t-1} + b_{\overleftarrow{h}}\right) \tag{9}$$

$$y_t = W_{\vec{h}y}\vec{h}_t + W_{\overleftarrow{h}y}\overleftarrow{h}_t + b_y \tag{10}$$

*Training* BRNN can be prepared with an indistinguishable algorithm from standard unidirectional RNN in light of the fact that there are no communications among the two kinds of state neurons and, in this manner, can be extended into a general encourage forward system. Be that as it may, if, for instance, any type of back-propagation through time (BPTT) is utilized, the forward and backward pass methodologies are marginally more muddled in light of the fact that the relic state and yielded neurons should never again be possible each one in turn. On the off chance that BPTT is utilized, the forward and reverse disregard the extended BRNN after some time nearly similarly consuming a general MLP. Some unique treatment is essential just toward the start and the finish of the preparation information. The forward state contributions and the regressive state contributions are not known. That could be made as a piece of the learning procedure; however, they are set arbitrarily with a fixed value (0.5). What is more, the neighborhood state subordinates for the forward states and for the regressive states are not known and are set here to zero, accepting that the data past that point are not critical for the present update, which is, for the limits, positive for the case. The preparation methodology for the unfolded bidirectional system after some time can be condensed as what takes after.

1. Forward Pass.
   Run all information for one at a time cut through the BRNN and decide all anticipated yields.

   - Do this forward pass only for forward states (from to) and in reverse states (from to).
   - Do this forward go for yield neurons.

2. Backward Pass.
   Ascertain the piece of the target work subsidiary for the time cut utilized as a part of the forward pass.

   - Do in reverse go for yield neurons.
   - Do in reverse pass only for forward states (from to) and in reverse states (from to).

3. Update Weights.
   In view of the above methods, the feature vectors are characterized by digits and words through speech signal.

# 4 Experimental results and discussion

In this segment, the execution of BRNN-SOM has been assessed and additionally contrasted execution along and existing algorithms such as RNN and DNN-HMM [36]. The execution is assessed with respect to SNR, MSE, and classification accuracy. The investigations were led utilizing Tamil queries taken from Forum for Information Retrieval and Evaluation (FIRE) dataset 2011. Fifty short Tamil title point queries uttered by 20 people with three reiterations aggregate of 3000 sentences were utilized for preparing, and 10 people with 2 redundancies aggregate of 1000 sentences were utilized for testing.

The principal metric utilized amid assessment of preprocessing algorithms is signal-to-noise ratio (SNR). SNR is utilized to measure how much a signal has been contaminated by noise. It is characterized as the ratio of signal power to the noise control ruining the signal. The SNR is figured in two ways: One is Pre-SNR, and other is Post-SNR which are acquired previously, then after the fact enforcing the preprocessing operation. De-noising is effective if Post-SNR is higher than Pre-SNR. Equation (11) shows the recipe used to gauge SNR.

$$SNR_{db} = 10 \log_{10}\left(\frac{P_{signal,db}}{P_{noise,db}}\right) = P_{signal,db} - P_{noise,db} \quad (11)$$

Mean square error (MSE) is utilized quantify the variations among esteems implied and the true being estimated. The MSE is determined with the help of Eq. (12),

$$MSE = \frac{1}{N}\sum_i (X_i - Y_i)^2 \quad (12)$$

where $x_i$ is the original signal, $y_i$ is the noisy signal and $x_i$ is estimated $x_i$ (noisy signal passed by means of de-noising algorithm). Lower MSE represents a closer match among the two signals.

The accuracy is computed with the help of Eq. (13). A high-accuracy value represents maximized speech recognition performance.

$$Accuracy(\%) = \frac{No.\ of\ words\ are\ correctly\ recognized}{Total\ No.\ of\ words} \times 100 \quad (13)$$

The noisy speech signals were improved by various speech pre-handling algorithms such as Gaussian filtering (GF) [18], hard and soft combined thresholding (HSCT) conspire [24], and proposed SGF plot. Three kinds of noise evacuation are centered, like, white noise, babble noise, and external noise. Three sorts of noise were evacuated by utilizing proposed SGF alongside existing two plans.

## 4.1 SNR comparison among various preprocessing schemes

The suggested SGF preprocessing scheme is distinguished with the current HSCT and GF methods with respect to final SNR for three sorts of noise removal, which are shown in Figs. 5, 6, and 7. The speech signals were utilized in this work, which were considered from the FIRE database. In x-axis an initial SNR is considered and y-axis final SNR is considered. It can be demonstrated that the sug-
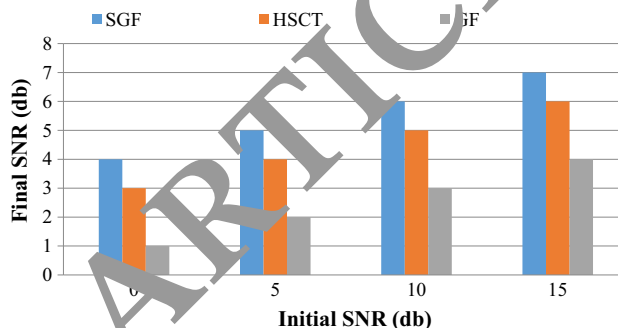


**Fig.** White-noise-removal-based SNR performance comparison among various preprocessing schemes
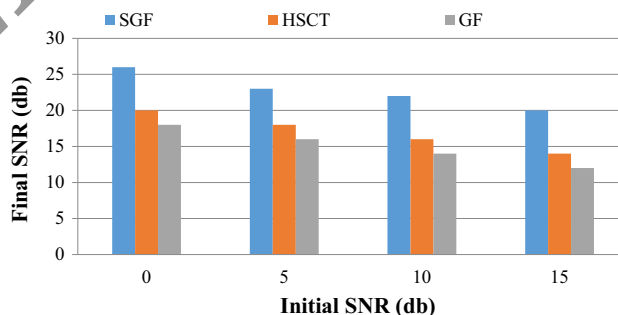


**Fig. 6** Babble-noise-removal-based SNR performance comparison among various preprocessing schemes
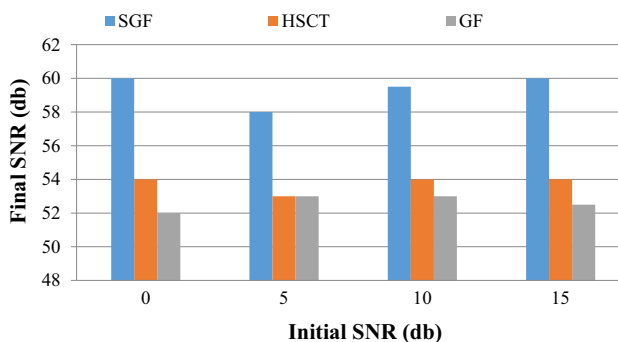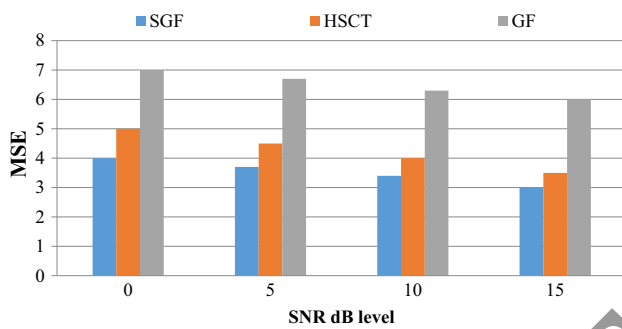


**Fig. 7** External-noise-removal-based SNR performance comparison among various preprocessing schemes

gested SGF approach accomplishes high final SNR value when distinguished with the other current speech enhancement methods.

## 4.2 MSE comparison among various preprocessing schemes

The suggested SGF preprocessing scheme is distinguished with the current HSCT and GF methods with respect to final MSE for three types of noise removal, which are shown in Figs. 8, 9, and 10. The speech signals are utilized here, and it is considered from the FIRE database. In $x$-axis an initial SNR is considered and $y$-axis final MSE is



**Fig. 8** White-noise-removal-based MSE performance comparison among various preprocessing schemes
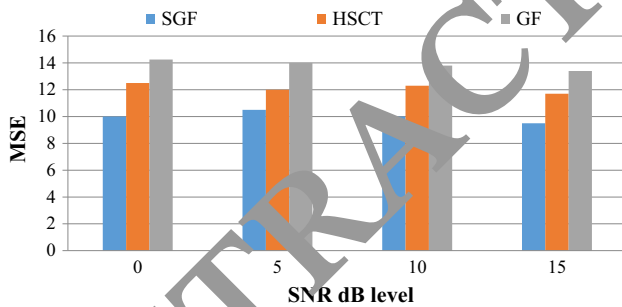


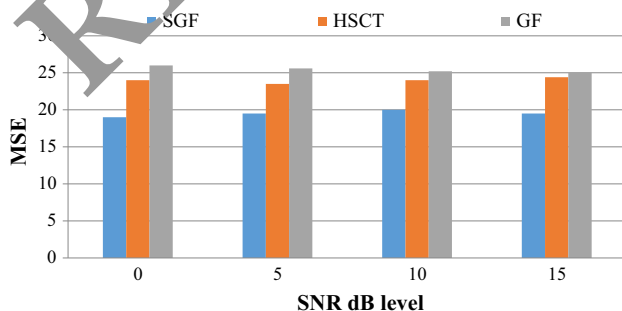**Fig. 9** Babble-noise-removal-based MSE performance comparison among various preprocessing schemes



**Fig. 10** External-noise-removal-based MSE performance comparison among various preprocessing schemes

considered. It can be proved that the suggested SGF approach accomplishes less final MSE value when distinguished with the rest of the current speech enhancement methods.

## 4.3 SNR comparison among various classification schemes

The suggested optimized BRNN-SOM is distinguished with the current RNN and DNN-HMM methods with respect to final SNR which is shown in Fig. 11. The speech signals are utilized here, which are considered from the FIRE database. In x-axis an initial SNR is considered and y-axis final SNR is considered. The SNR measure considers both residual noise level and speech degradation. It can be said that the proposed BRNN-SOM approach accomplishes high final SNR value when distinguished with the rest of the current speech enhancement methods. Due to the effectual preprocessing and feature extraction, the proposed scheme acquired better results.

## 4.4 MSE comparison among various classification schemes

The graphical indication of MSE performance comparison between the suggested and current algorithms is shown in Fig. 12. The speech signals utilized here were considered from the FIRE database. In x-axis an SNR level is considered and y-axis MSE is considered. The MSE measure



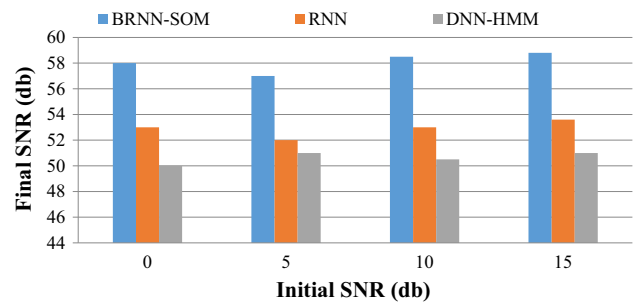**Fig. 11** SNR comparison among all ASR classification methods
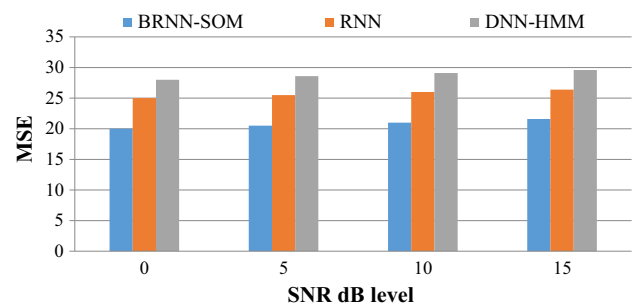


**Fig. 12** MSE comparison among all ASR classification methods

considers both SNR level and speech degradation. It proves the MSE performance of proposed BRNN-SOM scheme acquired less value when distinguished with the current RNN and DNN-HMM. Due to the effectual process of SOM, the proposed scheme acquires less error rate.

## 4.5 Accuracy comparison among various classification schemes

The graphical representation of accuracy performance comparison between the suggested and current algorithms is shown in Fig. 13. It proves the accuracy performance of proposed BRNN-SOM scheme acquires high-accuracy value of 93.6% when distinguished with the current RNN and DNN-HMM, and because of the effectual preprocessing and feature extraction, the proposed scheme acquires better results.

## 5 Conclusion

Lately, neural system has turned into an improved method for handling complex issues and dull assignments, for example, speech recognition. Speech is a characteristic and straightforward specialized strategy for individuals. Be that as it may, it is a to a great degree of mind-boggling and troublesome occupation to influence a PC to answer for the spoken commands. As of late, there is an earth-battering requirement for ASR framework to be created in Tamil and other Indian languages. In this paper, such a vital exertion is done for perceiving Tamil spoken words. To finish this assignment, feature extraction is done in the wake of utilizing required preprocessing systems. The most generally utilized PLP and MAR techniques are utilized to extricate the critical feature vectors from the upgraded speech signal, and they are provided as the contribution to the BRNN. The received system is trained with these input and target vectors. The experimental analysis demonstrates that the proposed plot accomplished 93.6% of exactness and better SNR and less MSE distinguished with current plans such as RNN and DNN-HMM. In future, this preparatory trial will
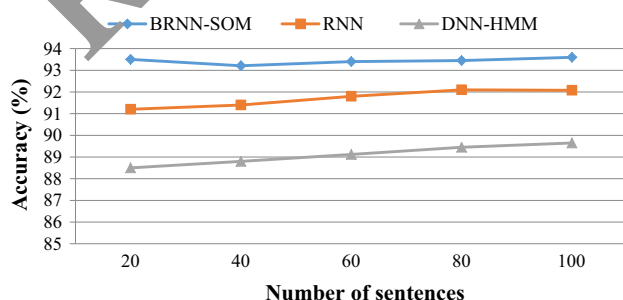

Fig. 13 Accuracy comparison among all ASR classification methods

creates ASR framework for Tamil language utilizing distinctive methodologies such as neural network based plans or with other cross-hybrid strategies.

## Compliance with ethical standards

**Conflict of interest** This statement is to certify that all authors have seen and approved the manuscript being submitted. We warrant that the article is the authors' original work. We warrant that the article has not received prior publication and is not under consideration for publication elsewhere. On behalf of all co-authors, the corresponding author shall bear full responsibility for the submission. The author(s) declare that there is no conflict of interest.

## References

1. Varatharajan R, Manogaran G, Priyan MK, Sundarasekar R (2017) Wearable sensor devices for early detection of Alzheimer disease using dynamic time warping algorithm. Cluster Comput. https://doi.org/10.1007/s10586-017-0977-2
2. Varatharajan R, Manogaran G, Priyan MK, Balaş VE, Barna C (2017) Visual analysis of geospatial habitat suitability model based on inverse distance weighting with paired comparison analysis. Multimedia Tools Appl. https://doi.org/10.1007/s11042-017-4768-9
3. Balan EV, Priyan MK, Gokulnath C, Devi GU (2015) Fuzzy based intrusion detection systems in MANET. Procedia Comput Sci 50:109–114
4. Devi GU, Balan EV, Priyan MK, Gokulnath C (2015) Mutual authentication scheme for IoT application. Indian J Sci Technol 8(26). https://doi.org/10.17485/ijst/2015/v8i26/80996
5. Manogaran G, Varatharajan R, Priyan MK (2018) Hybrid recommendation system for heart disease diagnosis based on multiple kernel learning with adaptive neuro-fuzzy inference system. Multimedia Tools Appl 77(4):4379–4399
6. Priyan MK, Devi GU (2017) Energy efficient node selection algorithm based on node performance index and random waypoint mobility model in internet of vehicles. Cluster Comput. https://doi.org/10.1007/s10586-017-0998-x
7. Varatharajan R, Manogaran G, Priyan MK (2017) A big data classification approach using LDA with an enhanced SVM method for ECG signals in cloud computing. Multimedia Tools Appl. https://doi.org/10.1007/s11042-017-5318-1
8. Devi GU, Priyan MK, Balan EV, Nath CG, Chandrasekhar M (2015) Detection of DDoS attack using optimized hop count filtering technique. Indian J Sci Technol 8(26):1–6. https://doi.org/10.17485/ijst/2015/v8i26/83981
9. Gokulnath C, Priyan MK, Balan EV, Prabha KR, Jeyanthi R (2015) Preservation of privacy in data mining by using PCA based perturbation technique. In: 2015 international conference on smart technologies and management for computing, communication, controls, energy and materials (ICSTM). IEEE, pp 202–206
10. Thota C, Sudarasekhar R, Manogaran G, Varatharajan R, Priyan MK (2017) Centralized fog computing security platform for IoT and cloud in healthcare system. In: Krishna Prasad AV (ed) Exploring the convergence of big data and the internet of things. IGI Global, Hershey, pp 141–154
11. Kumar PM, Gandhi U, Varatharajan R, Manogaran G, Jidhesh R, Vadivel T (2017) Intelligent face recognition and navigation system using neural learning for smart security in Internet of

Things. Cluster Comput. https://doi.org/10.1007/s10586-017-1323-4

12. Manogaran G, Varatharajan R, Lopez D, Kumar PM, Sundarasekar R, Thota C (2017) A new architecture of Internet of Things and big data ecosystem for secured smart healthcare monitoring and alerting system. Future Gener Comput Syst 82:375–387

13. Kumar PM, Gandhi UD (2017) A novel three-tier Internet of Things architecture with machine learning algorithm for early detection of heart diseases. Comput Electr Eng 65:222–235

14. Radha V, Vimala C, Krishnaveni M (2012) Continuous speech recognition system for Tamil language using monophone-based hidden markov model. In: Proceedings of the second international conference on computational science, engineering and information technology. ACM, pp 227–231

15. Radha V, Vimala C, Krishnaveni M (2011) Isolated word recognition system for Tamil spoken language using back propagation neural network based on LPCC features. Comput Sci Eng 1(4):1–11

16. Patel I, Rao YS (2010) Speech recognition using HMM with MFCC: an analysis using frequency spectral decomposition technique. Signal Image Process Int J (SIPIJ) 1(2):101–110

17. Chandrasekar M, Ponnavaikko M (2008) Tamil speech recognition: a complete model. Electron J Tech Acoust, article no. 20. http://www.ejta.org/en/chandrasekar2

18. Rojathai S, Venkatesulu M (2012) A novel speech recognition system for Tamil word recognition based on MFCC and FFBNN. Eur J Sci Res 85(4):578–590

19. Sigappi AN, Palanivel S (2012) Spoken word recognition strategy for Tamil language. Int J Comput Sci Issues 9(1):1694–0814

20. Sivaraj P, Rama M (2012) Recognition of isolated spoken words using DWT. Int J Eng Sci Res 2(9):1187–1196

21. Thangarajan R, Natarajan AM, Selvam M (2008) Word and triphone based approaches in continuous speech recognition for Tamil language. WSEAS Trans Signal Process 3(3):76–86

22. Saraswathi S, Geetha TV (2010) Design of language models at various phases of Tamil speech recognition system. Int J Eng Sci Technol 2(5):244–257

23. Karpagavalli S, Rani KU, Deepika R, Kousalya P (2012) Isolated Tamil digits speech recognition using vector quantization. Int J Eng Res Technol 1(4):1–12

24. Iswarya P, Radha V (2012) Speech based query processing architecture for Tamil-English in cross language text retrieval system. Int J Emerg Trends Eng Dev 7(2):437–442

25. Schafer R (2011) What is a Savitzky-Golay filter? IEEE Signal Process Mag 28:111–117 (lecture notes)

26. Savitzky A, Golay MJE (1964) Smoothing and differentiation of data by simplified least squares procedures. Anal Chem 36:1627–1639

27. Neumaier A, Schneider T (2001) Estimation of parameters and eigenmodes of multivariate autoregressive models. ACM Trans Math Softw (TOMS) 27(1):27–57

28. Lütkepohl H (2005) New introduction to multiple time series analysis. Springer, Berlin

29. Box GE, Jenkins GM, Reinsel GC, Ljung GM (2015) Time series analysis: forecasting and control. Wiley, Hoboken

30. Misra H (2006) Multi-stream processing for noise robust speech recognition. Doctoral thesis, Swiss Federal Institute of Technology (EPFL), Lausanne, Switzerland, March 2006

31. Chen R, Jamieson LH (1996) Experiments on the implementation of recurrent neural networks for speech phone recognition. In: Proceedings of the thirtieth annual Asilomar conference on signals, systems and computers, Pacific Grove, California, November, pp 779–782

32. Lee ST, Kim KC, Moon H, Cho JW (1991) Application of fully neural networks for speech recognition. In: Korea Advanced Institute of Science and Technology, Korea, pp 77–80

33. He J, Liu L (1999) Speaker verification performance and the length of test sentence. In: Proceedings on ICASSP 1999, vol 1, pp 305–308

34. Gingras F, Bengio Y (1998) Handling asynchronous or missing data with recurrent networks. Int J Comput Intell Organ 1(3):154–163

35. Schuster M, Paliwal KK (1997) Bidirectional recurrent neural networks. IEEE Trans Signal Process 45:2673–2681

36. Fredes J, Novoa J, King S, Stern RM, Yoma NB (2017) Locally normalized filter banks applied to deep neural-network-based robust speech recognition. IEEE Signal Process Lett 24(4):377–381