# Punjabi Speech Syllable Segmentation Using Vowel Onset Point Identification

**S. Kaur[1*], M.K. Gill [2]**

[1,2] Department of Computer Science and Engineering, Guru Nanak Dev Engineering College, Punjab, India

*Corresponding Author: shelly.ldh12@gmail.com, Tel.: 9465180728*

**Abstract—** Speech Recognition has been a wide area of research for a long time now. Researchers have been putting a lot of efforts and devised different methods for the same. For Speech Recognition system, speech signal is divided or segmented into some acoustic units like phonemes, syllables and word which will reduces the search space for unwanted signal or noise. This research work aims at developing an Automatic Speech Segmentation algorithm for Punjabi language which segments the signal into syllabes. For Automatic Speech Syllable Segmentation, a proposed technique detects the syllable boundaries using gamma tone filter and oscillator. In this proposed technique, valley picking picks the valley of the signal and gives the onset of the speech signal. Results of proposed technique was compared with the existing method which takes less time. After that Automatic Speech Classification algorithm classifies the signal into two classes either native or non native. For this, system had been trained using Artificial Neural Network (ANN) for estimating the parameter of Native and Non-Native spekers using Mel Frequency Cepstrum Coefficients (MFCCs) for feature extraction. The whole work was performed in Matlab2016a and the results generated as output with high accuracy.

*Keywords—* MFCC, ANN, MATLAB, Punjabi language, gamma tone fiter bank and oscillator.

## I. INTRODUCTION

Speech recognition technique has been explored widely for the past four decades as the degree of acceptance for such systems is high. Speech is considered to be one of the easiest and comfortable means of communication. It is an efficient way to recognize a person on the basis of speech. It is an important biometric authentication process [1]. The development workflow for speech recognition follows the acquisition of speech, feature extraction module, the recognition model and the testing module. In this work, implementation is done using MATLAB. MFCCs are used for feature extraction because it is designed using the knowledge of human auditory system and is used in every state of speech recognition system or art speech. It is used mostly as it is believed to mimic the behavior of human ear [2] [3]. Automatic Speech Classification is done in this work. The task may seem easy but for machines it can be challenging task due to high acoustic similarities among certain group of letters. Automatic Speech Segmentation into its acoustic units like phonemes, syllable and word is a difficult task. It is done vary carefully. There are various methods for Automatic Speech Segmentation those are Wavelet method, Artificial Neural Networks, Short term Energy, Word Chopper technique and Hidden Markov

Model which gives better results then the Manual Segmentation. Manual Segmentation takes more time and error prone.The speech recognition area aims to develop a technique and system which is used for voice input to machine based on advanced static voice modeling, automatic speech recognition is now widely used in tasks requiring human machine interfaces such as automatic call processing. If we see in past, computer scientists have been investigating ways and means to enable computers to understand and interpret human speech.

*Types of Speech*
1. Isolated Word: Isolated word recognizes that each utterance usually needs to be quiet on both sides of the sample windows. It accepts one word or one utterance at a time. It has "listen and not listen." Isolated utterance of this class could be better [4].
2. Connected Word: The connected word system is similar to isolated words, but allows separate words to take a minimum break between them together.
3. Continuous Speech: Continuous speech recognizers enable the person to speak almost naturally, whereas the computer determines the continuity. Recognition system with continuous speech abilities is one of the hardest to create

because it uses a special method to determine the pronouncement limit [5].

4.Spontaneous speech: At a fundamental level, it can be considered as natural trying to sound and not rehearsed. An ASR scheme with completely spontaneous speech capability should be able to manage a variety of natural voice features, such as words that are used together.

The main objective of this work is to identify the syllable boundary of the speech signal and then segmented into its smaller units which is the syllable using onset of the signal. For this work gamma tone filter bank and oscillator is used which picks the valley of the signal and gives the onset values and segmented the signal into syllables. After that, for identifying whether the word is spoken by the native or non-native person Classification technique is used. For the purpose, Artificial Neural Network classifier is used.

Rest of the paper is organized as follows, Section II contains the Related work, Section III contains Problem formulation, Section IV contains Methodology, Section V contains Results and Discussion of the proposed work and Section VI concludes the entire work and its future scope.

## II. RELATED WORK

**Geetha et al.** [10] explains segmentation of speech signals into linguistic units has enormous applications in the fields of recognition, synthesis, labeling and transcription and coding. In the proposed work, Tamil speech segmentation task has been carried out to identify the boundaries of syllable in the given speech utterance. A novel syllable segmentation method is proposed which uses the VOPs as anchor points to identify the position of consonant vowel unit and spectral difference are used find the syllable boundary.

**Mary et al.** [11] proposes an algorithm which evaluated in terms of detection accuracy, missed syllable and spurious syllables for different speech modes such as read, lecture and conversation mode of speech in Malayalam and Bengali.

**Ozseven et al.** [13] proposed a novel tool for speech feature extraction and classification. In this work, they developed a toolbox called SPAC for MATLAB based speech processing and feature extraction. There are various toolboxes prepared for this purpose in the literature, but they have advantages and weaknesses compared to each other.

**Sakran et al.** [6] shows that among the techniques investigated the short term energy technique which suffered from the problem of thresholding. The Word Chopper based segmentation cannot be used for segmentation of all words.

**Nasereddin et al.** [12] made a comparative study between classification techniques from ASR point of view, as well as, the translation approaches from MT point of view. Speech processing is considered to be one of the most important application area of digital signal processing process and performance.

## III. PROBLEM FORMULATION

In Punjabi, there are seven types of syllables. The pronunciation of Punjabi syllables spoken by native and non-native speakers have major differences. They pronounce same word in different way containing variable stress on vowels and consonants. In order to process Punjabi speech at syllable level, waveform of speech is to be segmented into smaller acoustic units. Ambiguity arises when same words are spoken by native and non-native speakers of the language [7].

Automatic speech segmentation is important for speech recognition because it reduces the search space effectively in automatic speech recognition [8]. The existing work has considered three different language samples namely Hindi, Odia and Bengali and their algorithm has not been tested for Punjabi language. Also, there algorithm considers three syllable forms namely V, CV and CCV [9]. But this does not considers the instances where more than one vowels occur in tandem. This can lead to unwanted or ineffective outcomes when such systems are realised for real world speech recognition systems. Thus, the speech recognition system needs to developed by modifying the existing one such that it can also consider the other syllables also which may lead to improve the accuracy of the system.

## IV. METHODOLOGY

In this work, Automatic Speech Segmentation is done on isolated Punjabi words using gamma tone filter bank and oscillator based segmentation for dividing the speech signal into syllable like units and the wav files are used for segmentation are 16 bit samples. Working of segmentation method is discussed below:

1. Input wav file is recorded at sampling frequency at 16 KHz from both native and non-native persons.
2. After recording the speech, signal preprocessing is done to extract the useful part of speech.
3. Features from isolated words of Punjabi speech signal were extracted using MFCC (Mel Frequency Cepstrum Coefficient). MFCC is the best feature extraction technique.
4. Gamma tone filter bank and oscillator based segmentation method segmented the speech signal into syllables and gives the desired output having syllable boundaries detected automatically.
5. At the end, Classification is done using Artificial Neural Network classifier for identify that the word is either spoken by the native or non-native person.
6. Figure 1 is the flow diagram of isolated Punjabi words Automatic Speech Segmentation and Classification steps involved during the methodology.
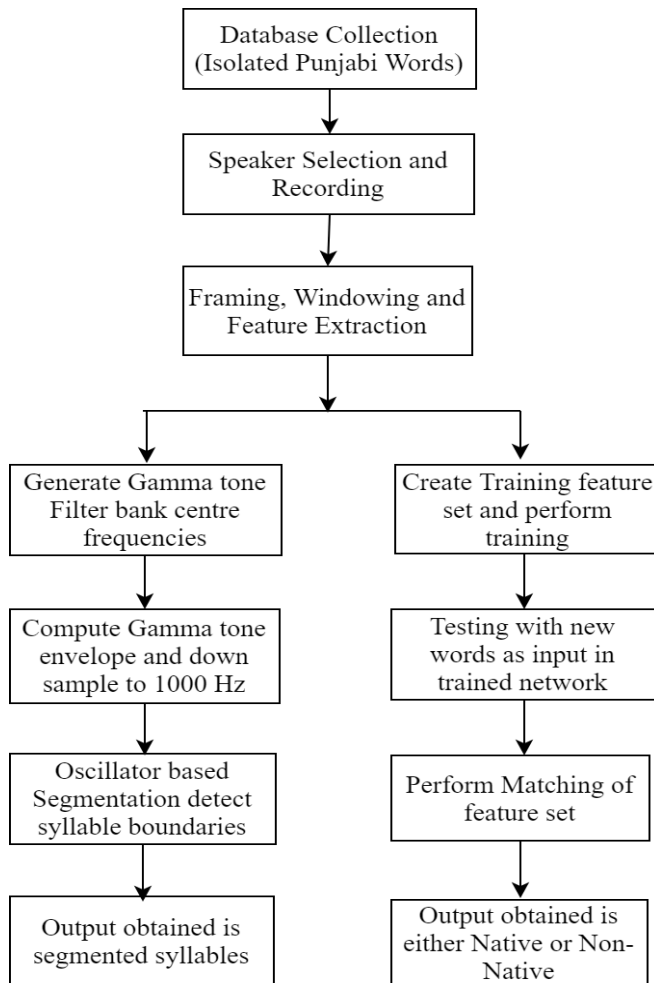
```
        ┌─────────────────────────┐
        │   Database Collection   │
        │ (Isolated Punjabi Words)│
        └─────────────────────────┘
                    │
        ┌─────────────────────────┐
        │  Speaker Selection and  │
        │        Recording        │
        └─────────────────────────┘
                    │
        ┌─────────────────────────┐
        │ Framing, Windowing and  │
        │   Feature Extraction    │
        └─────────────────────────┘
             │              │
   ┌──────────────────┐ ┌──────────────────┐
   │ Generate Gamma   │ │ Create Training  │
   │ tone Filter bank │ │ feature set and  │
   │ centre frequencies│ │ perform training │
   └──────────────────┘ └──────────────────┘
             │              │
   ┌──────────────────┐ ┌──────────────────┐
   │ Compute Gamma    │ │ Testing with new │
   │ tone envelope and│ │ words as input in│
   │ down sample to   │ │ trained network  │
   │ 1000 Hz          │ │                  │
   └──────────────────┘ └──────────────────┘
             │              │
   ┌──────────────────┐ ┌──────────────────┐
   │ Oscillator based │ │ Perform Matching │
   │ Segmentation     │ │ of feature set   │
   │ detect syllable  │ │                  │
   │ boundaries       │ │                  │
   └──────────────────┘ └──────────────────┘
             │              │
   ┌──────────────────┐ ┌──────────────────┐
   │ Output obtained  │ │ Output obtained  │
   │ is segmented     │ │ is either Native │
   │ syllables        │ │ or Non-Native    │
   └──────────────────┘ └──────────────────┘
```

Figure 1: Automatic Speech Segmentation and Classification steps

**Algorithm for Proposed Technique:**
**Step 1:** Collect Data in form of isolated words corpus in Punjabi language.
**Step 2:** Select Speaker along with recording.
**Step 3:** Preprocess the Signal
(a) Perform framing by dividing speech into time frames.
(b) Do windowing by using hamming window to eliminate discontinuities.
**Step 4:** Extract features using MFCC.
**Step 5:** Identify the syllable boundaries using Gamma tone filter bank and oscillator-based segmentation.

*A. Data Collection*
Punjabi is one of the language that used by most of the Punjab people speakes and some other states people also. Punjabi is a syllabic language contains seven types of syllables. In this thesis, data should be taken is isolated words corpus in Punjabi language. The word selected in such a manner that each word is a combination of vowels and consonants, which can be easily segmented into syllable boundaries [10, 11]. In

Punjabi language, there are seven types of syllables V, VC, CV, VCC, CVC, CCVC and CVCC where V and C represents vowel and consonant respectively. The word corpus we taken contains 150 isolated Punjabi words spoken by each individual persons either native or non native. Each word contains one or more than one syllable. Database contains 6000 words spoken by 20 native and 20 non native persons.

*B. Speaker Selection and Recording*
Multilingual speakers have been selecting for recording from others state. Speech sample was recorded from both native and non-native speakers. Native speakers are those whose mother tongue language is Punjabi and non-native speakers are from other state. Native and non-native speakers are males within the age limit of fifteen to thirty years. Speech sample contains 20 native male speakers whose mother tongue is Punjabi and 20 non- native male speakers. Data are recorded with the help of a unidirectional microphone using a recording tool audacity in a normal room with minimum external noise. The sampling rate used for recording is 16 kHz. All files are recorded using good quality microphone to make it free from noise. Software used recording purpose is Audacity and the recorded files are wave files within the time duration 0 to 3 sec.

*C. Signal Preprocessing*
It is very crucial to Pre-Process the Speech Signal in the applications where silence or background noise is completely undesirable. The speech is first divided into time frames consisting of an arbitrary number of samples. By introducing overlapping of the frames, the transition from frame to frame is smoothed. The speech samples are segmenting into small frame size. The number of samples used for each frame is 256. To calculate the number of frames, total numbers of samples in the input voice file are divided by 128. All time frames are then windowed with a Hamming window to eliminate discontinuities at the edges for subsequent Fourier Transforms [12] [13]. Generally Speaking, the use of short frame duration and overlapping frames is chosen to capture the rapid dynamics of the spectrum. Speech parameters are extracted on a frame- by-frame basis and the amount of overlap determines how quickly parameters can change from frame to frame. Window duration determines the amount of averaging used in power or energy calculation.

- *Function used for Framing in MATLAB*
Fs= 16 KHz, Fs is the sampling rate at which wav file is recorded.

$$\text{Frame duration} = fr\_length / Fs$$
$$= 320/16000$$
$$= 0.2 \text{ sec} = 20$$

- *Framing can be done as:*
        nbFrames = Ceil ((length(x)-N)/M);
        Frames = zeros (nbFrames+1, N);
Where, nbFrames is the number of samples per frame.

- *Function used for Windowing in MATLAB:*
                w = hamming (nbSamples);

### D. Feature Extraction

Feature Extraction is the basic requirement for all of the speech processing applications, and the main purpose is to derive descriptive attributes from the signal [12]. Feature extraction is a process of obtaining different features such as power, pitch and vocal tract configuration from the speech signal. In this work, MFCCs is used for feature extraction. Mel. Frequency Cepstral Coefficients (MFCC) technique is the robust and dynamic technique for speech feature extraction [14] [15].

- *Functions used for MFCC:*
1. FFTs i.e. Fast Fourier Transform function of Matlab is used to calculate power spectrum and generate Mel filter.
2. Convert linear scale frequency into mel scale frequency using:
                melmax = 2595 * log10 (1+ fmax/700);
3. Apply mel filters to Power Spectrum coefficients:
                melPowSpecs = w * PowSpecs;
4. Calculate MFCC using DCT (Discrete Cosine Transform) function.
                melCeps = dct(log(melPowSpecs));

### E. Automatic Speech Syllable Segmentation

In Automatic speech syllable segmentation, speech signal is segmented into syllable like units automatically [15]. In this thesis work, syllable segmentation is done automatically using proposed method. Gamma tone filter bank and oscillator-based segmentation is done to find the syllable boundaries of isolated Punjabi words.

`Gammatone Filter Bank:` Gammatone Filter Bank decomposes a signal by passing it through a bank of gammatone filters equally spaced on the ERB scale. Gammatone filter banks were designed to model the human auditory system.

- *Function used for gamma tone filter bank:*
    [bm, env, instp, instf] = gammatone_c (x, fs, cf, hrect)
Where x = input signal, Fs = sampling frequency, cf = centre frequency of the filter (Hz), hrect = half- wave rectifier, bm = basilar membrane displacement env = instantaneous envelope, instp = instantaneous phase (unwrapped radian), instf = instantaneous frequency (Hz)

Oscillator Based Segmentation
The detector is based on an interconnected network of leaky integrate-and-fire (LIF) neurons. Its principles are based on findings on the role of slow neural oscillations in auditory cortex for natural speech parsing. In essence, the network is composed of $n_E = 10$ excitatory and $n_I = 10$ inhibitory neurons.

- *Function used for Oscillator is:*
  [bounds,  bounds_  t,  env,  nuclei,  outs]  = thetaOscillator(ENVELOPE, f, Q, thr, verbose)

Inputs: ENVELOPE = N * F matrix, where N is the number of samples and F is the number of frequency bands, f = center frequency of the oscillator (Default = 5 Hz), Q = Q-Value of the oscillator (Default = 0.5), thr = detection threshold of the syllable boundaries (Default = 0.025).
Outputs:  bounds = cell array of boundaries (in 1000 Hz samples).
bounds_t = cell array of boundaries in time (seconds),env = cell array of oscillator envelopes,  nucleii = cell array of nuclei location (in samples).

### F. Automatic Speech Classification

For classification, artificial neural network is used, which defines whether speaker is native or non-native. Several neural network models including competitive learning and counter propagation are developed to identify individuals as either native or non-native speakers based on their accents [2]. Some important speech features, such as pitch period and the first three formant frequencies, are used as inputs to the neural networks [14] [17]. It is a feasible approach for an assisting automatic speech classification system in an environment in which different accents may be used [18].

In order to design and implement an Automatic Speech Classification system Neural Network is used. In MATLAB there was a Neural Network Toolbox that provided many of the functions needed to implement any type of neural network.

The model of neural net is trained on the training dataset using a supervised learning method i.e. gradient descent with momentum and adaptive learning rate back propagation. If the features of test data are matched with the training data the desired results are obtained. When matching is performed testing features of the data sets was matched with the training data.

## V. RESULTS AND DISCUSSION

Speech database contains 150 Punjabi words from each native and non-native speakers. The Speech signal was preprocessed first to remove the unwanted data or noise from the speech samples and to extract the feature of the speech. In this thesis work, syllable segmentation of a speech signal is done using MATLAB version 2016a.  The word selected in such a manner that each word is a combination of vowels and consonants. Speech sample contains 20 native male speakers whose mother tongue is Punjabi and 20 non-native male speakers.

### A. Results of Framing and Windowing

The speech is first divided into time frames consisting of an

arbitrary number of samples. The speech is first divided into time frames consisting of an arbitrary number of samples. By introducing overlapping of the frames, the transition from frame to frame smoothed. All time frames are then windowed with a Hamming window to eliminate discontinuities at the edges for subsequent Fourier Transforms.
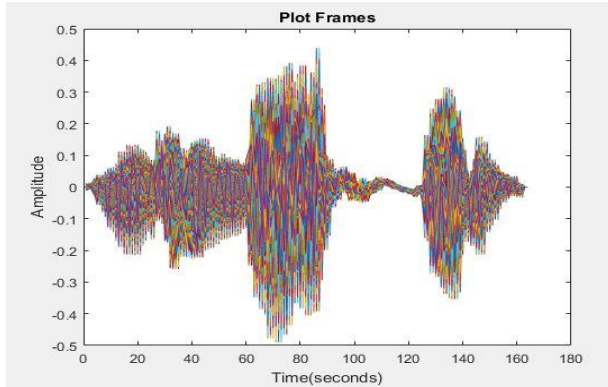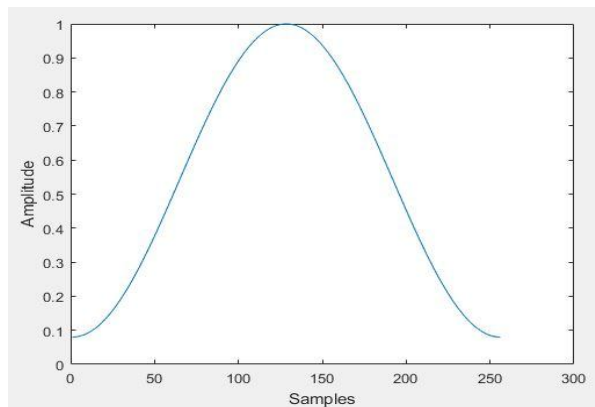


Figure 2: Framing of the word 'ਗੁਰਦਾਸਪੁਰ'



Figure 3: Hamming Window of the word 'ਗੁਰਦਾਸਪੁਰ'

### B. Results of Feature Extraction

Mel-frequency Cepstral Coefficient (MFCC) technique is often used to extract important feature of sound file. Two feature matrices are formed after implementing MFCC technique. The output of featured matrix is given below.



Figure 4. Hamming Window of the word 'ਗੁਰਦਾਸਪੁਰ'

### C. Segmentation using Spectrogram and Segmentation parameters

In the Spectrogram and parameters-based segmentation, Spectrogram of the speech signal was plotted first using the values of parameters i.e. window, no. of overlaps, NFFT and MINDB. Using the Segmentation parameters like its cutoff, maximum magnitude in the spectrogram and frequency segmentation is done. The red color markers in the signal depicts the segmentation of the signal, which could not clearly detect the boundaries of the syllables.
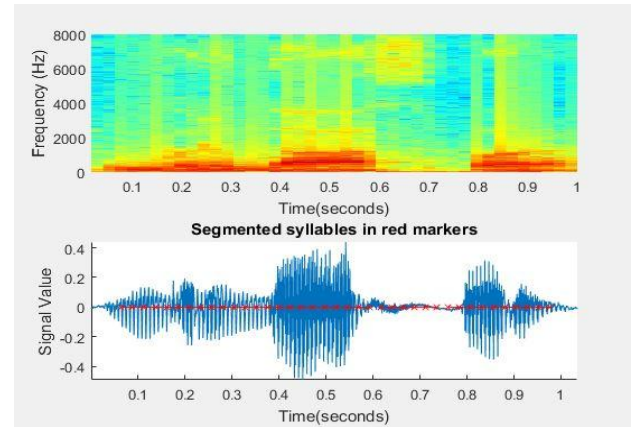


Figure 5: Segmentation using spectrogram and

segmentation parameters of the word 'ਗੁਰਦਾਸਪੁਰ'

### D. Automatic Speech Syllable Segmentation using gammatone filterbank and oscillator

In Automatic Speech Syllable segmentation, vertical lines made from the oscillator gives the syllable boundaries of the speech signal and the red color umbrella depicts the number of segments in the speech signal. In the automatic segmentation of the word 'ਗੁਰਦਾਸਪੁਰ' is segmented into three parts forming the syllables of same type CVC. The red color umbrella depicts that the word has three segments.
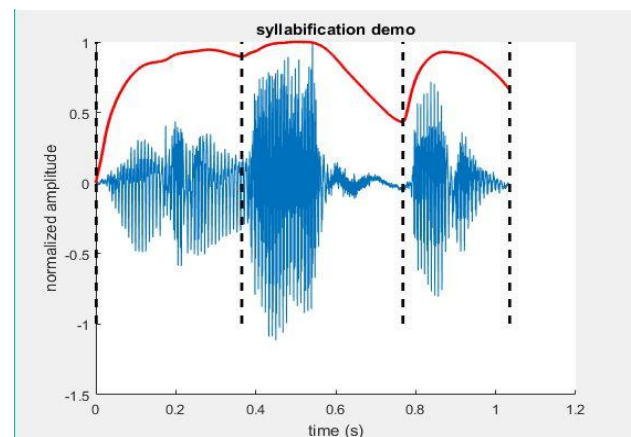


Figure 6: Automatic speech syllable segmentation of

the word 'ਗੁਰਦਾਸਪੁਰ'

### E. AUTOMATIC SPEECH CLASSIFICATION

For classification, artificial neural network is used, which defines whether speaker is native or non-native. Training feature set was created using Mel Frequency Cepstrum Coefficients of isolated Punjabi word spoken by 14 speakers. Other remaining speaker's words are used for the testing of network. When to test the data, features of testing set are extracted using the Mel Frequency Cepstrum Coefficients. The feature set of trained network performs matching with the testing set. If the features of the trained network and testing data set are matched, we get the desired results either it is native or non-native.
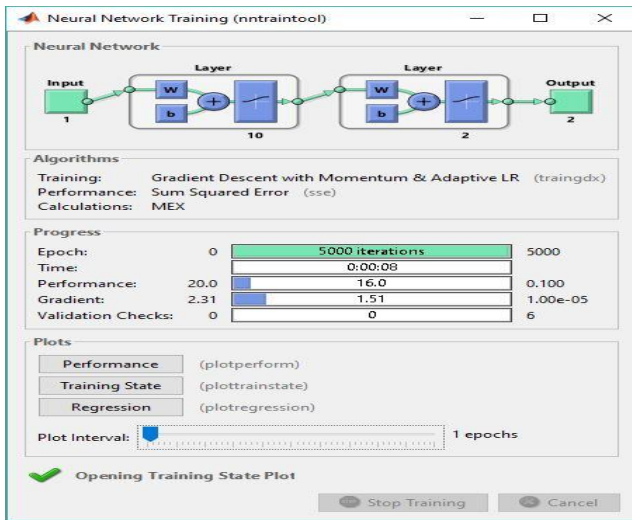


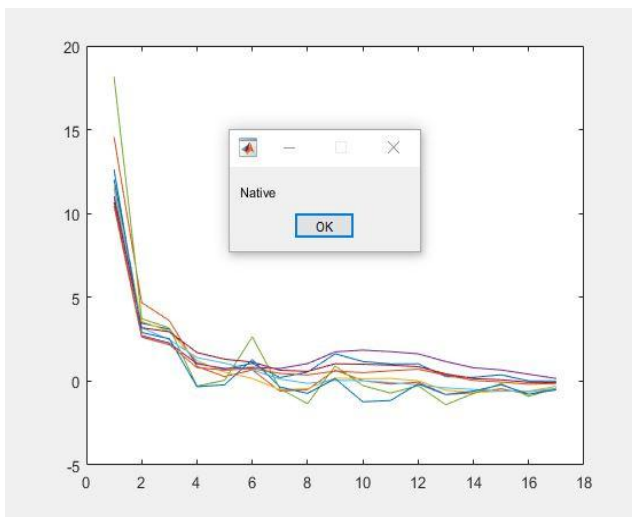Figure 7. ANN Classifier training toolbox



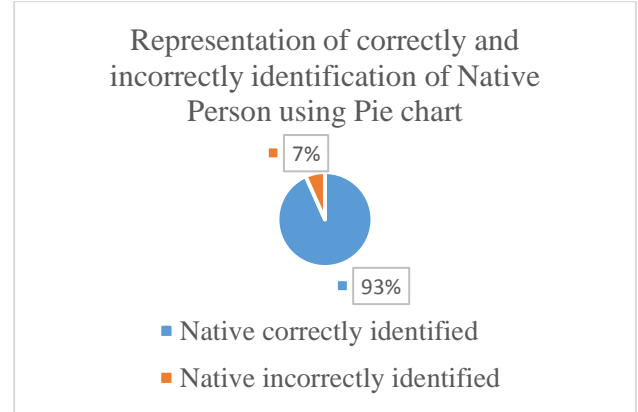Figure 8: Classifier identified as Native Speaker



Figure 9: Representation of correctly and incorrectly identification of Native Person using pie chart
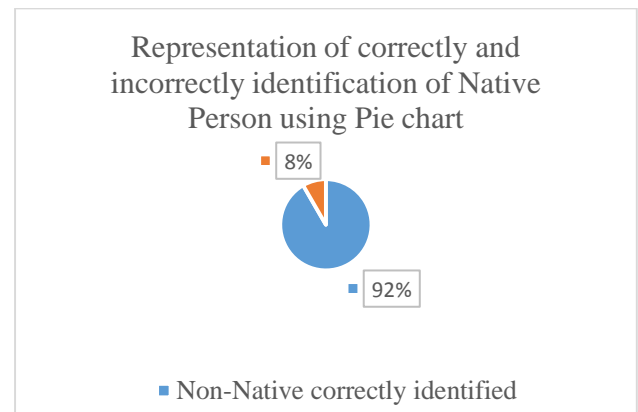


Figure 10: Representation of correctly and incorrectly identification of Non-Native Person using pie chart

### F. Comparison of Existing and Proposed Technique on the basis of Time
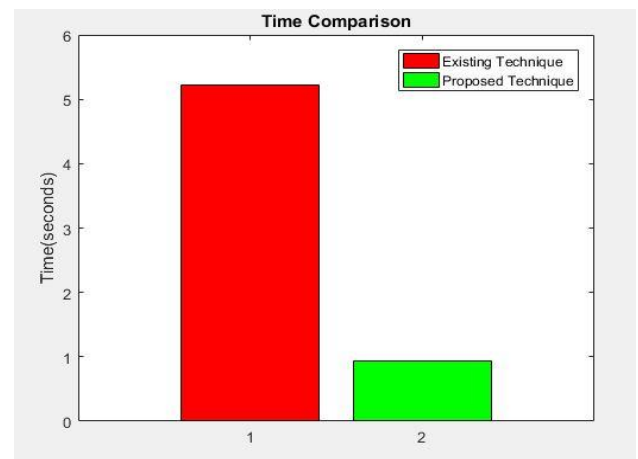


Figure 11: Comparison of Existing and Proposed technique on the basis of time

Table 1. Time Comparison between Existing and Proposed Technique

| S.No | Isolated Punjabi Words | Existing Technique (Time in (sec)) | Proposed Technique (Time in (sec)) |
|------|------------------------|-----------------------------------|-----------------------------------|
| 1 | ਪਾਕਿਸਤਾਨ | 6.1697 | 0.8810 |
| 2 | ਮੰਦਾ | 5.4427 | 0.7729 |
| 3 | ਗੁਰਦਾਸਪੁਰ | 5.2281 | 0.9345 |
| 4 | ਸਾਡਾ | 6.9582 | 0.9020 |

## VI.   CONCLUSION AND FUTURE SCOPE

Segmentation of speech signals into its acoustic units has huge applications in the fields of recognition, synthesis, labeling and coding. This research work proposed an algorithm for finding syllable boundaries using gamma tone filter bank and oscillator-based segmentation. The Speech samples were taken from both native and non-native speakers. In this thesis, data collected is isolated words corpus in Punjabi language. The syllable boundaries are detected automatically by using gamma tone filter bank and oscillator. The vertical lines made from oscillators represents accurate boundary of syllables. The red color umbrella depicts number of segments of syllables. Experimental results clearly show that the results given by the proposed technique are more efficient as compared to the existing technique. Further V, CV, VC and CVC syllables are more clearly identified as segmentation of these are gives more correct results as that of syllables CCVC and CVCC, which provide incorrect results because consonants get separated. The existing technique was time consuming and error prone while the proposed method takes less time and gives us more accurate and desired results. For this wok, only Punjabi language is used whereas multiple languages in future can also be used for automatic speech syllable segmentation. The database collection takes only isolated words however; in future, continuous speech would be used to enhance this work further. The wok can also be extended for the syllable containing consonant-consonant part which sometime gets separated in our proposed method. Artificial Neural Network classifier was used for the classification of native and non native person by extracting its features for training and testing data and gives accuracy 93.1%. The accuracy can also be improved using other parameters for the training and the testing. In future, other feature extraction techniques like Linear Predictive Coding, Linear Predictive Cepstral Coefficients, and Wavelet Packet Decomposition can also be used for classification.

## REFERENCES

[1]   Y. Youhao," Research on Speech Recognition Technology and Its Application,"in the proceedings of the 2012 International Conference on Computer Science and Electronics Engineering Research, vol. 6, no. 12, pp. 306-309, 2012.

[2]   K. Amino, T. Osanai, *"Native vs. non-native accent identification using Japanese spoken telephone numbers,"* Speech Communication, vol. 56, no. 1, pp. 70–81, 2014.

[3]   M. Wester, C. Mayo, "Accent rating by native and non-native listeners," in the proceedings of the 2014 ICASSP IEEE International Conference Acoustically Speech Signal Processing, no. i, pp. 7699–7703, 2014.

[4]   D. B. Hanchate, M. Nalawade, M. Pawar, V. Pophale, P. K. Maurya,*"Vocal Digit Recognition using Artificial Neural Network,"* IEEE Journal, vol. 7, no. 10, pp. 88–91, 2010.

[5]   L. Bouafif and K. Ouni, "A speech tool software for signal processing applications," in the proceedings of the 2012 6th International Conference Science Electronics Technology Information Telecommunication SETIT, pp. 788–791, 2012.

[6]   E. Sakran, S. M. Abdou, S. E. Hamid, M. Rashwan, *"A Review : Automatic Speech Segmentation,"* International Journal of Computer Science and Mobile Computing, vol. 6, no. 4, pp. 308–315, 2017.

[7]   P. Kumari, D. Shakina Deiv, M. Bhattacharya, "Automatic speech recognition of accented Hindi data," in the proceedings of the 2014 International Conference on Computation of Power, Energy, Information and Communication(ICCPEIC), pp. 68–76, 2014.

[8]   A. Kaur, E. T. Singh, *"Segmentation of Continuous Punjabi Speech Signal into Syllables,"* World Congress on Engineering and Computer Science, vol. I, pp. 20–23, 2010.

[9]   S. P. Panda, A. K. Nayak, *"Automatic speech segmentation in syllable centric speech recognition system,"* International Journal of Speech Technology, vol. 19, no. 1, pp. 9–18, 2016.

[10]  K. Geetha, R. Vadivel, *"Syllable Segmentation of Tamil Speech Signals Using Vowel Onset Point and Spectral Transition Measure,"* Automatic Control and Computer Sciences, vol. 52, no. 1, pp. 21–25, 2018.

[11]  L. Mary, A. P. Antony, *"Automatic syllabification of speech signal using short time energy and vowel onset points,"* International Journal of Speech Technology, pp. 571– 579, 2018.

[12]  S. S. Tirumala, S. R. Shahmiri, A. S. Garhwal, *"Speaker Identification feature extraction methods: A Systematic Review",* International Journal of Elsevier, Vol(90),pp. 250-271, 2017.

[13]  T. Ozseven, M. Dugenci, *"SPeech ACoustic (SPAC): A novel tool for speech feature extraction and classification",* International Journal of Elsevier, 136, pp. 1-8, 2018.

[14]  M. R. Gamit, K. Dhameliya, Dr. N. S. Bhatt, *"Classification Techniques for Speech Recognition",* International Journal of Emerging Technology and Advanced Engineering, vol. 5, no. 2, pp. 58-63, 2015.

[15]  C. P. Bharat, A. A. Desai, *"Segmentation of Gujarati words from Continuous spoken Gujarati Speech Signal,"* VNSGU Journal of Science and Technology, vol. 4, no. 1, pp. 106-112, 2015.

[16]  B. Barhate, D. Sisodiya, R. Deore, *"Applications of Speech Recognition: For Programming Languages,"* International Journal of Scientific Research in Computer Science and Engineering, vol. 6, no. 1, pp. 6-8, 2018.

[17]  Madan, D. Gupta, *"Speech Feature Extraction and Classification,"* International Journal of Computer Applications, vol. 2, no. 1, pp. 10-15, 2014.

[18]  H. K. Soni, *"Machine Learning – A New Paradigm of AI,"* International Journal Scientific Research in Network Security and Communication, vol. 7, no. 3, pp. 31- 32, 2019.

**Authors Profile**

*Ms. Sukhjinder Kaur* studied Bachelor of Technology in Computer science and Engineering from Guru Nanak Dev engineering College Ludhiana (Punjab) in 2017. She is currently pursuing Master of Technology in Computer Science & engineering from Guru Nanak Dev Engineering College, Ludhiana (Punjab). Her main research focuses on segmentation of speech signal into its acoustic units automatically and classification of that signal automatically.

*Ms. Manjot Kaur Gill* received a Bachelor degree in Computer Science Engineering in 2006 and a Master degree in Computer Science & Engineering in 2009. She is with Guru Nanak Dev Engineering College, Ludhiana (Punjab) India as Assistant Professor, in CSE department since 2008-2019. Her research area is Speech Recognition and Natural Language Processing.