

Automatic Recognition and Detection of Words in Speech- A Review

Kavya K S

P.G. Student, School of Computer Sciences M G University, Kottayam, Kerala India

ABSTRACT: Speech is the primary and the most convenient means of communication between people. Automatic Speech Recognition provides a way for natural communication between human and machine. It has become an interesting and challenging area. It helps to provide a capability to a machine for responding properly to spoken language. The main aim of this study is to understand the various approaches used for the recognition and detection of speech so that we can find out the methods that gives better accuracy and performance.

KEYWORDS: Automatic Speech Recognition, DTW algorithm, HMM, Keyword Spotting, MFCC .

I. INTRODUCTION

Speech is the vocalized form of communication used by humans, which is based upon the syntactic combination of items drawn from the lexicon. Each spoken word is created out of the phonetic combination of a limited set of vowel and consonant speech sound units (phonemes). These vocabularies, the syntax that structures them, and their sets of speech sound units differ, creating many thousands of different, and mutually unintelligible, human languages. Speech is achieved by compression of the lung volume causing air flow which may be made audible if set into vibration by the activity of the larynx. This sound can then be made into speech by various modifications of the supralaryngeal vocal tract.

Speech recognition (SR) is the inter-disciplinary sub-field of computational linguistics that develops methodologies and technologies that enables the recognition and translation of spoken language into text by computers. It is also known as automatic speech recognition (ASR). ASR is technology that allows a computer to identify the words that a person speaks into a microphone or telephone and convert it to written text [9]. The aim of ASR research is to allow a computer to recognize speech in real-time as human understands [2]. The words that are spoken by any person independent of vocabulary size, noise, speaker characteristics and channel conditions is a major challenge in ASR. The accuracy rate of speech recognition depends on the feature extraction techniques as well as the training methods used. A speech interface would support many valuable applications, i.e., telephonedirectory assistance, spoken database querying for novice users, hands busy applications in medicine, officedictation devices, or automatic voice translation into foreign languages [8]. ASR is technology that allows a computer to identify the words that a person speaks into a microphone or telephone and convert it to written text

Keyword spotting has become an important branch of audiomining. Keyword spotting is well suited to data mining tasks, process large amount of speech such as real time monitoring and to audio document indexing. The task of locating the occurrences of a given keywords K in a speech utterance is termed as keyword spotting (KWS). It plays an important role in audio indexing and speech data mining applications. The method is applied mainly for a large amount of spoken documents must be searched to learn whether they contain some specific words. The fast detection of the words and information about the exact location eliminate a lot of human work in such tasks like audio data mining [3]. Hence, Keyword spotting is well suited to data mining tasks for process large amount of speech because keyword spotting requires significantly less processing power than transcription, and can therefore run at considerably faster speeds.

International Journal of Innovative Research in Science, Engineering and Technology

(A High Impact Factor & UGC Approved Journal)

Website: www.ijirset.com

Vol. 6, Issue 8, August 2017

This Paper is organized as follows. Section II describes. The major researches works carried out in the automatic speech recognition and the word spotting, Finally, Section III presents conclusion of the review done in this work.

II. RELATED WORK

1. Chandra, E. [3] in this paper focused on Keyword Spotting System for Tamil Isolated Words using Multidimensional MFCC and DTW Algorithm. Tamil month names are recorded in noisy environment by a female speaker. Each month name in Tamil has been repeated 25 times so that totally 300 words are recorded using AUDACITY speech software. A preprocessing is made for not only noise reduction, but also normalization. Then 12D, 26D and 39D MFCCs are calculated for all utterances and are stored as templates. Finally, DTW is used as a pattern matching algorithm due to its speed and efficiency in detecting similar patterns. Experiments have been conducted with multiple MFCCs to achieve best word identifier. The system is designed for Tamil language but all its modules except the spoken word references developed are language independent. This paper shows detailed clarification about audio mining and the keyword spotting. Audio mining also called audio searching technique is used to search audio files for occurrences of spoken words or phrases. The major advantage of audio mining is the ability to process and search audio data thousands of times faster for large files than a human could - because the human would have to listen to each word of the audio. There are two main phases in audio mining system: training and template matching. During the training phase, a training vector is generated from the speech signal of each word spoken by the user. The training vectors extract the spectral features for distinguishing different classes of words. Each training vector can serve as a template for a single word or a word class. These feature vectors are stored in a database for subsequent use in the template matching phase. During the keyword matching phase, the user speaks any word that is to be identified in the trained templates. Feature vector is generated for that word and compared with the trained templates using pattern matching algorithm.

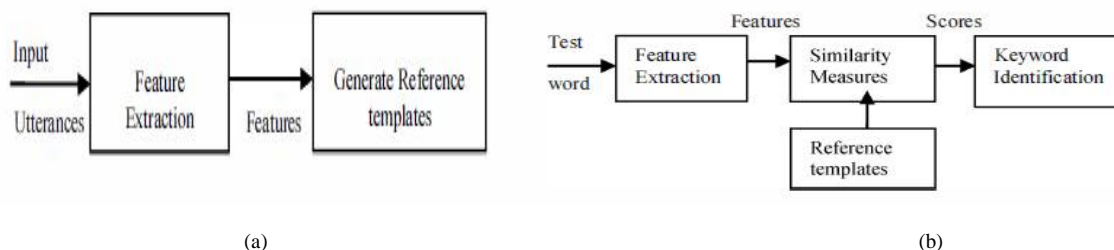


Fig 1. Block Diagram Of Audio Mining System (a) Feature Extraction Phase (b) Feature Matching Phase

The preferred technique for feature extraction is MFCC wherein the features are generated by transforming the signal into frequency domain. In general, cepstral features are more compact, discriminable, and most importantly, nearly decorrelated and therefore, they can provide higher baseline performance over filter bank features. In phase 2 DTW is used. Dynamic time warping is an efficient algorithm optimizes to find the nonlinear alignment path between two sequence. their distance is obtained by time scaling one of the signals nonlinearly so that it aligns with the other. For keyword spotting, a prototype of the keyword is stored as template and compared to each of word in the incoming speech utterance. The experimental results shows that the number of occurrences of given keywords are identified and maximum of 89.7% average accuracy of keyword spotting is obtained with 39 dimensional MFCC. The 12, 26 and 39 dimensional MFCC are used for keyword spotting and it is found that 39 dimensional MFCC are more efficient and has greater accuracy than 13 and 26 dimensional MFCCs.

2.K. Anoop, P. Vivek and V. L. Lajish [1] in this work propose a keyword spotting approach for content based video indexing and retrieval system for classifying and indexing broadcast news videos in Malayalam. here used auditory modality to identify keywords to analyze the content of the news videos. A template based keyword spotting approach is then proposed for the purpose. Log Energy Filter Coefficient (LEFC) feature is the major component in the keyword spotting. The audio signals are segmented into 10ms frames and applied the filter bank for extraction of eight dimensional LEFC feature. In feature matching stage, the Dynamic Time Warping (DTW) of the template and the test audio data feature vector is performed. The similarity distance measures are computed and it is normalized by the

International Journal of Innovative Research in Science, Engineering and Technology

(A High Impact Factor & UGC Approved Journal)

Website: www.ijirset.com

Vol. 6, Issue 8, August 2017

sum of energy in each template. Best matches from all templates are stored according to the distance measure and finally the desired keyword will be spotted based on this distance measure.

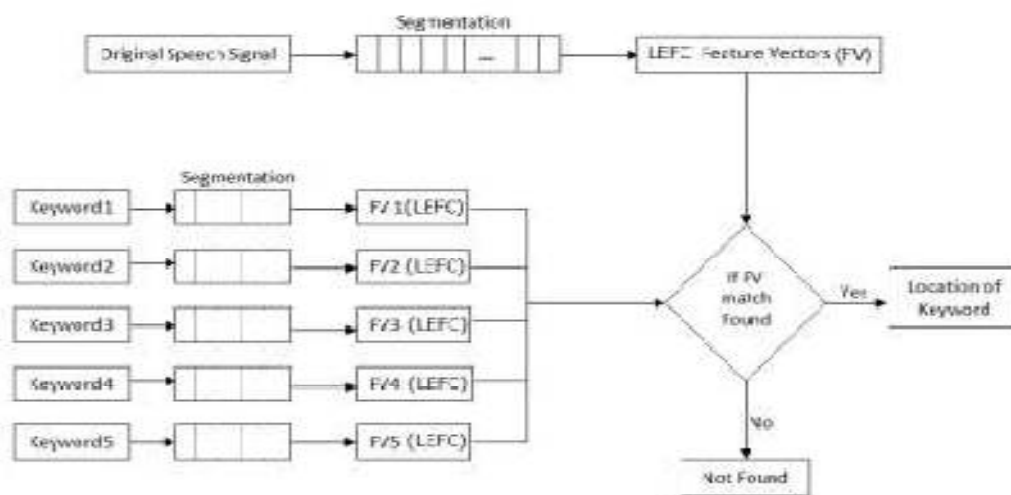


Fig 2. Keyword spotting using template matching

The final results shows that, template based keyword spotting using the LEFC feature can be effectively used for the speaker depended keyword spotting, while it will be less effective with the speaker independent keyword spotting.

3. Kurian, Cini, and Kannan Balakrishnan.[4] This paper presents a report on the development of a speaker independent, continuous transcription system for Malayalam. The system employs Hidden Markov Model (HMM) for acoustic modeling and Mel Frequency Cepstral Coefficient (MFCC) for feature extraction. It is trained with 21 male and female speakers in the age group ranging from 20 to 40 years.

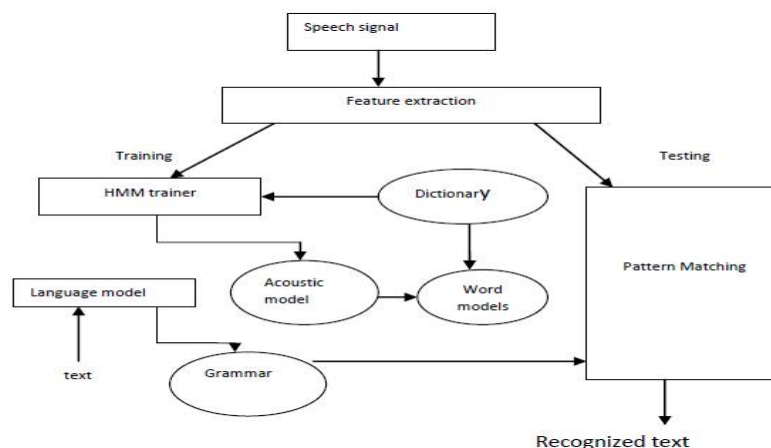


Fig 3. System Architecture of transcription system

During training phase the HMM trainer creates acoustic models from three components; feature vectors, language models, and dictionary. Word models which are built from language models and pronunciation dictionary are used

International Journal of Innovative Research in Science, Engineering and Technology

(A High Impact Factor & UGC Approved Journal)

Website: www.ijirset.com

Vol. 6, Issue 8, August 2017

for pattern matching in the testing phase along with the syntax and semantics of the language. The simulation results shows that the transcription system for Malayalam achieved while working with medium size vocabulary.

4. Satori, Hassan, et al. [6] in this developed Arabic voice recognition system using entirely Arabic environment. Both training and recognizing process use Arabic characters. Here both training and test were based on CMUSphinx4 system, which is HMM-based, speaker-independent, continuous recognition system capable of handling large vocabularies. This approach for modeling Arabic sounds in The CMUSphinx system consisted of generated and trained acoustic and language models with Arabic speech data. The dictionary adopted in the experiments was made using Arabic characters. The training process consists of: convert the audio data to a stream of feature vectors, convert the text into a sequence of linear triphone HMMs using the pronunciation dictionary, and find the best state sequence or state alignment through the sentence HMM for the corresponding feature vector sequence. For each senone, gather all the frames in the training corpus that mapped to that senone in the above step and build a suitable statistical model for the corresponding collection of feature vectors. The circularity in this training process is resolved using the iterative Baum-Welch or forward-backward training algorithm. The final results the recognition rate calculated that is found with better accuracy.

5. Radha, V. [5] This research work presents a speaker independent isolated speech recognition system for Tamil language. The most flexible and successful approach to speech recognition so far has been Hidden Markov Model (HMM) which is implemented in this research work. The HMM method provides a natural and highly reliable approach of recognizing speech for a broad range of applications. The system usually consists of two phases. They are called pre-processing and post-processing. Pre-processing involves feature extraction and the post-processing stage comprises of building an acoustic model, phonetic lexicon or the pronunciation dictionary and language modeling. Initially the speech waveform is processed to produce a new representation as a sequence of vectors containing values of features or parameters. The parameter values extracted from raw speech are used to build acoustic models. The Dictionary provides pronunciations for words found in the Language Model. The pronunciations break words into sequences of sub-word units found in the Acoustic Model. The other component is called a language model, which gives the probabilities of sequences of words. The pre-processing or front-end of a speech recognizer is responsible for the extraction of the features from the speech waveform. The post processing consists of modeling of the system. The performance of speech recognition system is generally measured in terms of word error rate, which is the ratio between misclassified words, and total number of tested words. The results shows that the obtained error rate for this research work is 0.88.

6. Moneykumar, Maya, et al [2] in this paper Malayalam Word Identification For Speech Recognition System proposed. Here also the MFCC feature extraction will be taken and the HMM model used for the final identification. Here the word identification system designed for Malayalam language uses a syllable based segmentation approach. Instead of training the system with independent words, we use syllables and phoneme combinations for training and identification. This eliminates the problem of maintaining large set of training data, as different words common syllable as well as phone segments. New words can be added to the vocabulary without building new models for the existing syllables and phonemes corresponding to the word. The experiment results shows that syllable based word identification system for Malayalam using HTK on the Linux platform was performed using MFCC and HMM. The training was conducted for 40 vocabularies of bi-syllable words. The system was trained to identify the utterances within the training vocabulary, with a moderate accuracy.

III. CONCLUSION

The review of Automatic Recognition and Detection of Words in Speech for different spoken language has been made and all the papers give different innovative ideas. Major works in the recognition and detection of speech carried out in Hidden Markov Model (HMM), this gives better performance compared to other existing model. And also the Mel Frequency Cepstral Coefficients (MFCC) can be used as the better choice for the feature extraction process. In

International Journal of Innovative Research in Science, Engineering and Technology

(A High Impact Factor & UGC Approved Journal)

Website: www.ijirset.com

Vol. 6, Issue 8, August 2017

future works large vocabulary speech recognition will be considered for improving the accuracy of the recognition system.

REFERENCES

- [1] K. Anoop , P. Vivek and V. L. Lajish A Keyword Spotting Approach For ContentBased Indexing And Retrieval Of Malayalam News Videos, NSA-2015 Goa, India, October 7-9, 2015
- [2] Moneykumar, Maya, et al. "Isolated Word Recognition System for Malayalam using Machine Learning." *12th International Conference on Natural Language Processing*. 2015.
- [3] Chandra, E. "Keyword spotting system for Tamil isolated words using Multidimensional MFCC and DTW algorithm." *Communications and Signal Processing (ICCSP), 2015 International Conference on*. IEEE, 2015.
- [4] Kurian, Cini, and Kannan Balakrishnan. "Automated transcription system for malayalam language." *International Journal of Computer Applications* 19.5 (2011): 5-10.
- [5] Radha, V. "Speaker independent isolated speech recognition system for Tamil language using HMM." *Procedia Engineering* 30 (2012): 1097-1102.
- [6] Satori, Hassan, et al. "Investigation arabic speech recognition using CMU sphinx system." *Int. Arab J. Inf. Technol.* 6.2 (2009): 186-190.
- [7] Ghai, Wiqas, and Navdeep Singh. "Analysis of automatic speech recognition systems for indo-aryan languages: Punjabi a case study." *Int J Soft Comput Eng* 2.1 (2012): 379-385.
- [8] Kurian, Cini, and Kannan Balakrishnan. "Speech recognition of Malayalam numbers." *Nature & Biologically Inspired Computing, 2009. NaBIC 2009. World Congress on*. IEEE, 2009.
- [9] Iqbaldeep Kaur, Navneet Kaur, Amandeep Ummat, Jaspreet Kaur, Navjot Kaur "Automatic Speech Recognition: A Review" *IJCST* Vo 1. 7, Issue4, oct- dec2016.