

A peer reviewed international journal ISSN: 2457-0362

www.ijarst.in

A ROBUST DEPENDENT AND INDEPENDENT SPEECH RECOGNITION BY USING ASR SYSTEM

N. Naresh Babu¹, N. Keerthana², G. Prachi², K. Sai Kirthana², G. Harshitha²

¹Assistant Professor, ²UG Student, ^{1,2}Department of Electronics and Communication Engineering ^{1,2}Malla Reddy Engineering College for Women, Maisammaguda, Hyderabad, Telangana, India

Abstract: It is proposed in this research to combine the discrete wavelet transform (DWT) with to feature extraction in order to enhance speech/speaker A Spectral (ASR) detection. Relative Different components of speech, such as speech recognition, speaker identification and speaker synthesis, all are connected to this. The goal of this research is to investigate the use of HMM with voice recognition. This method aims to improve performance besides introducing more features from either the speech signal including using parallel computations, which results in an increase between the recognition rate as well as computational speed both for the clean as well as noisy speech signals again for proposed method compared towards the HMM model.

Keywords: HMM Model, DWT, median filter RASTA-PLP.

1. INTRODUCTION

The most difficult part of building a voice recognition system was getting the recognition rate as accurate as possible while yet spending as little time as feasible on training and validation the system. We were able to solve these issues by selecting a decent feature extraction approach and implementing vour neural network classifiers in parallel. Herein lay the nub of our project. Combining DWT with RASTA-PLP allows for feature extraction. multi-resolution. Their multi-scale analytical capabilities of both the DWT make it ideal for processing non-stationary data like voice. RASTA-key PLP's benefit is that they are resistant to noise. That method's goal is to improve the new method's performance with feature extraction, resulting in a greater detection accuracy than that achieved by combining DWT and MFCCs-based methods.

2. EXISTING SYSTEM

This same wavelet-based HMM model methodology is what we used to propose speech/speaker recognition throughout this present method. The voice signal is used to create the wavelet basis. Its primary purpose was to satisfy a time-andbandwidth-consumption product. Even before FFT, it must have been employed to do the preprocessing. In just this case, some noise has been added.

Voice information may be utilised to identify individual speaker based on the anatomy of both the vocal tract, which really is unique to each individual. Speaker recognition was the procedure of identifying a person based on sound they make. Voice falls within the area called biometric identification because of the inherent distinctions inside the speaker's physical structure. There are several benefits to using one's voice as both a form of identification. Remote person authentication is indeed a big benefit. A speaker recognition system will have the same training and testing steps as every other pattern recognition system. In order to get a handle on how the system processes speech, training is necessary. Recognizing is done by testing. Figure 2.1 depicts the training period as just a block diagram. For develop the reference models, feature vectors reflecting the



speaker's voice characteristics were retrieved using training utterances.

The degree to which the test utterance's related feature vectors match the others in the reference was determined utilising some matching approach during testing. The choice is made based on the level each match. This testing phase's block diagram can be seen here.



Fig.2.1: The block diagram of training phase.



Fig.2.2: The block diagram of the testing phase

Every speech recognition system's extraction & selection of such an acoustic parametric description is important to its overall performance and accuracy. These MFCC (mel-frequency cepstrum coefficients) On either a mel-frequency scale, when are the cosine transforms of either the real logarithm of this with the short-term energy spectrum defined.

3. PROBLEM IDENTIFICATION

- Recognition process is less compare to mimicry voice
- Noise is very high

4. PROPOSED SYSTEM

Speech recognition was one of the ideas we floated. This RASTA-PLP

(relative spectral perceptual linear prediction) approach is being used in the method.High pass filters throughout the filter bank output may be used by Rasta to sluggish eliminate the variations throughout the components of either the mg and mg issues caused by communication channels.



Fig4.1: System flow chart

4.1. Discrete Wavelet Transform

When it comes to the DWT, two separate scales and places are picked depending on their respective capabilities. Powers of 2 are used to rescale and dilate original mother wavelet, whereas integers are used to translate it. You might demonstrate f(t), which specifies the space

> A peer reviewed international journal ISSN: 2457-0362

www.ijarst.in

JARST

- of square integrable functions, to be f(t)
- \in L2(R), by using the notation

$$f(t) = \sum_{j=1}^{L} \sum_{k=-\infty}^{\infty} d(j,k) \Psi(2^{-j}t-k) + \sum_{\substack{k=-\infty\\ -k}}^{\infty} a(L,K) \emptyset(2^{-L}t)$$

This function $(t)\phi$ is referred to as the scaling function, whereas the function $\phi(t)$ was referred to as one of the mother wavelet. Functions as a group

$$\{\sqrt{2^{-1}}\phi(2^{-j}t-l)|j\leq L,j,k,L\in Z|\}$$

Using Z as an orthonormal basis of L2, we get L2 (R). For a scale L, the approximation coefficients were denoted by a(L,k), whereas the detail coefficients have been designated by a(J,k). The following is an expression for the approximations and detail coefficients:

$$d(L,K) = \frac{1}{\sqrt{2^L}} \int_{-\infty}^{\infty} f(t) \, \emptyset \left(2^{-L}t - k \right) dt \dots (2)$$
$$d(j,k) = \frac{1}{\sqrt{2^j}} \int_{-\infty}^{\infty} f(t) \, \Psi \left(2^{-j}t - k \right) dt \dots (3)$$

These coefficients above may well be explained by the projection fl(t) on f(t)gives perfect calculation (in the sense of least error energy) of f(t). These coefficients a(L, k) might be used to generate this projection.

$$f_1(t) = \sum_{k=-\infty}^{\infty} a(L,K) \emptyset(2^{-L}t - k) \dots (4)$$

Scale 1 lowers, and the approximations becomes finer, eventually reaching f(t) when 1 drops to zero. It is quite possible to explain the difference in approximation among scales of approximation, fl+1(t) and scales of

approximation, fl(t), through using coefficients d(j,k).

$$f_{l+1}(t) - f_1(t) = \sum_{k=-\infty}^{\infty} d(l,k) \varphi (2^{-l}t - k) \dots (5)$$

Using those relations, given $a(L, k) \& \{d(j, k) \mid j \leq L\}$, It is obvious that now the approximation could be constructed at any size. As a result, the wavelet transform produces a rough approximation of the original signal.fL(t) (given a(L, k)) and a number of layers of detail {fj+1(t)-fj(t)| j < L} (given by {d(j, k) | j \leq L}). The closer the actual scale is to the next finer layer of information, the more accurate the final result becomes.

4.2. Vanishing Moments

This number of vanishing moments in either a wavelet demonstrates the smoothness and flatness of just a wavelet function like a frequency response, respectively (filters used to compute the DWT). P-vanishing-moment wavelets satisfy the following formula.

$$\int_{-\infty}^{\infty} t^m \varphi(t) dt = 0 \quad for = 0, \dots, p-1$$

Or equivalently,

$$\sum_{k} (-1)^{k} k^{m} c(k) = 0 \text{ for } m$$
$$= 0, \dots, p-1$$

The more vanishing moments that were in the representation given smooth signals, the faster the wavelet coefficients decayed. A wavelet with a high number of point sets effectively represents signals, which is beneficial for coding purposes. As the number of vanishing moments as well as the size of something like the wavelet filters rise, so does the computational difficulty of calculating the DWT coefficients. In the image below, we can see the low-frequency approximation

A peer reviewed international journal ISSN: 2457-0362

www.ijarst.in

coefficients cA1 as well as the high-frequency detail coefficients cD1.

IJARST



Fig.4.2: Filtering operation of the DWT

4.3. Implementation Using Filters

Fast Wavelet Transform method implementation is shown in the following formula:

 $\emptyset(t) = \sum_k c(k) \varphi(2t - k)$(6)

$$\varphi(t) = \sum_{k} (-1)^{k} c(1-k) \phi(2t-k)....(7)$$
$$\sum_{K} C_{K} C_{K-2M} = 2\delta_{0,m} \dots \dots \dots \dots (8)$$

The scaling function was defined by the twin-scale relation (also known as the dilation equation). Using the scaling function φ , the wavelet is expressed. In order therefore for wavelet to really be orthogonal to both the scaling function as well as its translating function, this third equation must be satisfied.

4.4. Multilevel Decomposition

In iterations, successive approximations may be deconstructed in turn to break down a single signal into several lower resolution components. This decomposition process called iterative. This same wavelet decomposition tree is indeed the name given to this structure.



Fig4.3: Decomposition of DWT co-efficient

This is the wavelet decomposition of both the signals for level j: [CAj, CDj,..., CD1]. Information may be discovered by looking at such a signals wavelet decomposition tree.

4.5 Signal Reconstruction

This technique may be used to recreate or create a new signal from the old one (IDWT). Approximation coefficients cAj, CDj, are used to approximate CAj-1 before up sampling & filtering with both the reconstruction filters are used to rebuild CAj-1.



Fig4.4:Filterlevels

This aliasing caused inside this wavelet decomposition step is countered by the design of both the reconstruction filters. Low and high-pass decomposition filters

A peer reviewed international journal ISSN: 2457-0362

www.ijarst.in

(Lo_ R and Hi_ R) are part of a quadrature mirror filters system (QMF).

IJARST



Fig4.5:Quadrature mirror filters (QMF)

4.6. Feature extraction: Wavelet Energy

The approximation as well as the detail each maintain a specific amount of energy whenever the signal was decomposed that use the wavelet decomposition technique. The wavelet accounting vector as well as the wavelet decomposition vector may be used to retrieve this energy. It is a ratio since it compares between original and decomposed signals, which yields the computed energy. This can only be discovered by a great deal of trial and error. Level 2 coefficients comprise the majority of the speech signal's associated data, therefore they are used to derive the wavelet decomposition coefficients.

5. RESULTS

This section gives the detailed analysis of simulation results. Further, the performance of proposed ASR system is compared with the state of art approaches using same dataset.

5.1 Dataset

The dataset contains the voice samples from five different persons. From each

person 9 samples are collected with multiple phrases. Totally, 45 phrases are collected in entire dataset. Further, 80% of dataset is used for training, 20% of dataset is used for testing the ASR system.

5.2 Subjective Performance

Figure 5.1 shows the original audio with 2.5 seconds of time. Figure 5.2 shows the frequency of original audio, which is having the frequency of -3k to +3k. Figure 5.3 shows the spectrogram of MFCC. Figure 5.4 shows the spectral features of MFCC, which contain 8th order feature, cepstral feature, log-mfcc features. Figure 5.5 shows the recognized person details.



Figure 5.1. Original audio



Figure 5.2. Frequency spectrum.

IJARST

A peer reviewed international journal ISSN: 2457-0362

www.ijarst.in



Figure 5.3. Spectrogram (a) original speech, (b) MFCC.



Figure 5.4. Spectral features of MFCC.



Figure 5.5. Recognized outcome.

5.3 Objective evaluation

Table 5.1. Performance comparison.

Metric	HDN	CNN	DBN	Propos
	N [1]	[2]	[7]	ed
				ASR
Accurac	92.55	94.7	94.2	99.98
y (%)	0	50	00	
Sensitiv	90.30	91.2	95.7	99.973
ity (%)	7	67	15	
Specific	94.77	94.1	95.6	99.83
ity (%)	0	04	69	
F-	93.46	92.8	95.2	99.67
measur	6	25	57	
e (%)				
Precisio	93.26	94.4	90.6	99.29
n (%)	2	28	88	
MCC	92.94	95.2	90.5	99.45
(%)	7	41	83	
Dice	90.76	90.6	92.7	99.37
(%)	2	24	27	
Jaccard	93.56	94.3	91.6	99.46
(%)	5	64	14	



Figure 5.6. Performance comparison.

Table 5.1 shows the performance of proposed ASR system with conventional approaches such as HDNN [1], CNN [2], DBN [7]. The proposed DLCNN based ASR resulted in superior performance for all metrics. Figure 5.6 shows the graphical representation of Table 5.1. Figure 5.7 and Figure 5.8 shows the graphical representation of GUIS.



Fig5.7: Voice Registration GUI



Fig 5.8: Voice Verification GUI Applications

- 1)Education System
- 2) Industrial
- 3) ATM and Banking
- 4)Electronic (Laptop,Mobile)

6. CONCLUSION

Methods for pre-processing the voice signal vary. The most accurate feature for word identification was found using wavelet transformations, as shown in the paper. **International Journal For Advanced Research**

In Science & Technology



A peer reviewed international journal ISSN: 2457-0362

www.ijarst.in

Respectively clean and noisy voice signals have shown this to be true. Wavelets are employed to test for dependent and independent variables in noisy speech, as well as the RASTA-PLP approach produces superior findings. verification tests were conducted and a suggested algorithm was found to have an accuracy rate of roughly 90% Comparison with the HMM model was unmistakable.

FUTURE SCOPE

Our future work concentrate on increasing the accuracy of the model and also to complete the system as a whole for the railway system we took dataset for. Our model works only for trained users and the future work is to extent by using neural networks is provide better security

REFERENCES

- [1]. SD.B.Paul,"Speech Recognition usingHidden Markov Model.International Magazine of Computers and Technology, No. 8, 2014 Vol. 13.
- [2].M.A.Anusuya , S.K.Katti"Speech Recognition by Machine: A Review" International journal of computer science and Information Security 2009.
- [3]. L.R.Rabiner and B.H.jaung ," Fundamentles of Speech Recognition Prentice-Hall, Englewood Cliff, New Jersy,1993.
- [4]. Mahdi Shaneh, and Azizollah Taheri,"Voice Command Recognition System Based on MFCC and VQ Algorithms", World Academy of Science, Engineering and Technology 57 2009
- [5] H. Combrinck and E. Botha, "On the mel-scaled cepstrum,"Video Communications, SPIE Electronic Imaging, UCRL Conf. San Jose, vol.200706, Jan 2004.
- [6] Electronic Engineering, University of Pretoria., Journal of Computer Science 3 (8): 608-616, 2007 ISSN 1549-3636.
- [7] Ahmad A. M. Abushariah, Teddy S.

Gunawan, Othman O. Khalifa"English Digits Speech Recognition System Based on Hidden Markov Models", International Islamic University Malaysia, International Conference on Computer and Communication Engineering (ICCCE 2010), 11-13 May 2010, Kuala Lumpur, Malaysia.

- [8] AnjaliBala,ABHIJEETKumar,Nidhika Birla,"Voice command recognition System Based on MFCC and DTW",International Journal of Engineering Science and Technology,Vol.2(12),2010
- [9] Ibrahim Patel,Dr.Y.Srinivasa Rao, , "Speech recognition using Hidden Markov Model With MFCCSubband Technique." 2010 International Conference on Recent Trends in Information,Telecommunication and Computing.
- [10] ShumailaIqbal,Tahira,Mehboob,Malik ,"Voice Recognition using HMM with MFCC for secure ATM",IJCS Vol.8,Issue 6 Nov 2011
- [11] Vimala C, Dr.V.Radha, "A Review on Speech Recognition Challenges and Approaches", World of Computer Science and Information Technology Journal (WCSIT) ISSN: 2221-0741 Vol. 2, No. 1, 1-7, 2012
- [12] Lawrence R. Rabiner, Fellow, IEEE
 'A Tutorial On Hidden Markov
 Model And Selected Applications
 In Speech Recognition,
 Proceedings Of TheIEEE, Vol. 77,
 No. 2, February 1989.