## Speaker Adaptation on Hidden Markov Model using MFCC and Rasta-PLP and Comparative Study

#### Shweta Bansal<sup>1\*</sup>, Atul Kumar<sup>2</sup> and S. S. Agrawal<sup>1</sup>

<sup>1</sup>KIIT College of Engineering, Gurgaon - 122102, Haryana, India; bansalshwe@gmail.com, ss\_agrawal@hotmail.com <sup>2</sup>Ansal University, Gurgaon - 122003, Haryana, India; atulkumar@ansaluniversity.edu.in

#### Abstract

This work compares the performance of the Mel-Frequency Cepstral Coefficient (MFCC) and Perceptual Linear Prediction (PLP) features for developing a text-dependent speaker identification system. Continuously spoken Hindi speech sentences have been used to train the HMM models using HTK toolkit for each speaker separately. The experiments have been performed using a set of 200 continuously spoken sentences with vocabulary of 20000 isolated words using a database of 100 speakers. The results show an accuracy of 92.26% recognition when PLP features have been used and accuracy of 91.18% for MFCC features. A confusion matrix has been created for all the 20 test speakers based on the recognition scores obtained for each of these speakers and their confusion with other speakers. Performance has been compared in the closed set and open set conditions of testing and as it is expected, the performance in the closed set condition is far better than the open set. We propose that if PLP features are used in place of MFCC, they may provide improvement in speaker than the open set. We propose that if PLP features are used in place of MFCC, they may provide improvement in speaker identification accuracy by reducing the cases of false acceptance.

Keywords: Hindi Speech, HMM, MFCC, RASTA-PLP, Speaker Identification

#### 1. Introduction

A speech signal not only contains information about text message but about the identity of the speaker. In speaker recognition, therefore, an attempt is made to extract the features and parameters from the signal which are helpful to identify or verify a speaker<sup>1</sup>. In the past, several experiments have been performed using MFCC features<sup>2</sup> but less number of experiments using PLP features<sup>3</sup>. These results have shown varying performances. Richard M. Stern et al. have found WER of 35.1% for MFCC and 38.0% for PLP features, for NRL Spine database<sup>4</sup>, but very few experiments have been done to compare the performance of speaker recognition using both MFCC and PLP features on the same database-particularly using Hindi database.

In the present paper, an attempt has been made to conduct speaker identification experiments from MFCC and PLP features. HTK toolkit<sup>5</sup> has been used to model the HMM classifier. The results obtained on using both the feature extraction techniques have been compared.

### 2. Overview of the System

Figure1 gives the overview of our approach for this study. The input speech is passed through front end processing unit. It includes the alignment of audio files and their transcription. Necessary processing is done to extract MFCC and PLP features from the speech signal. These feature vectors help to create the acoustic models containing all the feature vectors and duration of the words which exist in a particular sentence.

The lexicon module along with the language model plays the supportive role to find the system accuracy. The lexicon module contains all the lexicons (phones) contained in our database. Language model has been used to capture the possible variations in pronunciation uttered by a particular speaker, used to tied-states in HMM classifiers. We have tried to develop a phone-based trigram model to capture the sequence of phones in a word and then words in a particular sentence. With the help of all these inputs, the system follows the matching approach using Viterbi algorithm to find the best accuracy score of the words in a particular sentence<sup>6</sup>. The 1-best hypothesis of a sentence has been taken.





## 3. Methodology

#### 3.1 Database Collection

The present system uses a data-set of 200 phonetically rich Hindi sentences recorded by 100 native Hindi speakers (50 male and 50 female) within the age group of 18 to 30 years. Recording was carried out in a sound treated room environment having S/N > = 40 db. These were sampled at the rate of 16 bit-16 kHz. Eighty speakers were used for training the system and the other twenty speakers for testing the system in open-set condition. Similarly, the system has been trained for one hundred speakers and tested by using twenty speakers out of these 100 speakers in closed-set condition.

#### **3.2 Pronunciation Dictionary**

The following steps have been followed for the creation of pronunciation dictionary<sup>7</sup>:



#### 3.3 Transcription

For parallel alignment of audio and text data of various speakers, phonetic transcription has been done manually. For alignment purpose, all wave files were converted into text for the 100 speakers.

#### 3.4 Creation of Phone Sets

To capture the sound units in a word, phone-sets of Hindi have been created. 45 phones exist in our database, shown in Figure 2.

 $\mathfrak{M}[\mathbb{D}]$ ,  $\mathfrak{M}[\alpha]$ ,  $\mathfrak{P}[I]$ ,  $\mathfrak{F}[i]$ ,  $\mathfrak{I}[\circ]$   $\mathfrak{I}[\circ]$ ,  $\mathfrak{I}[$ 

 Table 3: Perception of Hindi consonants spoken by non native speakers and perceived by Hindi natives

### **3. ACOUSTIC ANALYSIS**

As the process of speech perception hinges on the acoustic spectrum, the formant frequencies provide enough information to distinguish the speech sounds. We have **Figure 2.** Hindi phone set.



Figure 2. Hindi phone set.

#### 3.5 Annotation at Phoneme Level

Speech database of 100 speakers were annotated at phoneme level using PRAAT software tool to capture the duration of individual phonemes. Figure 3 represents the samples of annotated files of two Hindi words (Himalaya and Vadiyaan) at phoneme level.





#### **3.6 Parameterization**

For extracting the relevant information from speech spectra MFCC and RASTA-PLP have been used for speaker diarization and compare their results using HTK<sup>8</sup>.

# 3.6.1 *Mel Frequency Cepstrum Coefficient* (*MFCC*)

Mel Frequency Cepstrum Coefficients are widely used for parameterization<sup>9</sup>. For MFCC extraction, the source kind is wave and the target kind is MFCC\_0<sup>10</sup>. The Hamming Window size used in the experiment has been taken to be 25 ms with a frame shift of 10 ms. Pre-emphasis co-efficients have been set as 0.97. Twenty-two band pass filters have been taken to capture all the co-efficients related to our database. For MFCC extraction, Logarithmic compression has been taken into consideration.

#### 3.6.2 Relative Spectral Transform-Perceptual Linear Prediction (RASTA-PLP)

RASTA-PLP is a technique of warping spectra to minimize the differences between speakers while preserving the important speech information<sup>11</sup>. It is an approach completely based on perceptual linear prediction.

The whole RASTA-PLP extraction is illustrated in the following figure:

For PLP extraction, the source format is wave and the target kind is PLP\_0<sup>11</sup>. Similar to MFCC, in this case also the Hamming window size of 25 ms and 22 bandpass filters have been taken to capture all the co-efficients from the database. For PLP extraction, cubic root compression have been taken into account.



Figure 4. Block diagram of RASTA speech processing technique.

#### 3.7 Speaker Modeling using HMM

We trained the recognition engine using MFCC and RASTA-PLP features for a database of 200 sentences for all the 100 speakers separately. For speaker identification, two models have been created - one is acoustic model which estimates the means and variances of a HMM in which each state output distribution is a single component Gaussian and the other is text-dependent language model<sup>12</sup> which is phone based trigram model to capture the sequence of phones in a word and then sequence of words in a particular sentence. To determine more accurate parameters, Baum-Welch re-estimation has been used<sup>13</sup>.

We perform single-pass retraining in which the state/ component occupation probabilities are calculated using an existing model and training set, but the new model parameters are calculated using a new training set<sup>14</sup>.

#### 3.8 Testing of the Recognition Engine

To recognize unknown test words, the likelihood of each model generating the word is calculated using Viterbi Algorithm<sup>15</sup> and the most likely model identifies the word. The testing is done separately with both MFCC and RASTA-PLP features.

## 4. Experimental Results

Figure 5 shows the performance of conducting the recognition tests. For example, it can be seen that if PLP features are used in place of MFCC, they may provide improvement in speaker identification accuracy and the cases of false acceptance may be reduced.



Figure 5. Comparison of MFCC and PLP feature.

The experiments were performed for all the hundred speakers in both closed and open-set. For closed-set comparison, the system is trained with 100 speakers, out of which 20 speakers have been used to test the system. For open-set comparison, the system is trained with 80 speakers and the remaining 20 speakers have been taken as a test data. Minimum threshold figure for accuracy has been considered as 60. Although the sentence contains varying number of words, the total number of words are 20000. The comparison with correct recognition and the false acceptance are shown in Table 1.

Test	Training	Correct	t	Correct Rejection				
Speakers	Utterances	Recogn	ition (%)		(%)			
		MFCC	RASTA-	MFCC	RASTA-			
			PLP		PLP			
20	20000	91.2	92.3	8.8	7.8			
(Closed								
Set )								
20	16000	64.2	74.6	35.8	25.4			
(Open								
Set)								

Table 1. Performance of system in closed and open set

The above results show that performance of system is significantly better when RASTA-PLP is used, particularly in case of open-set comparison. 

 Table 2(a). Confusion matrix of speaker identification

 engine (for MFCC features) in closed set condition



Table 2(b).Confusion matrix of speaker identificationengine (for PLP features) in closed set condition



#### 4.1 Confusion Matrix

In order to find out the individual performance of test speakers, confusion matrices has been made for the results obtained using MFCC and PLP features. These are shown in Table 2 (a and b) and Table 3 (a and b) for the closed set and open set experiments respectively. The confusion matrix show that in some cases the accuracy is higher for a test speaker even though the reference speaker is not same and providing ambiguity between themselves and the correct speaker<sup>16</sup>. For example – Speaker 5 gives 92.6% accuracy when tested with speaker 5 in closed set,

while speaker 50 gives 93.2% accuracy (Table 2(a)), which is more than the correct speaker.

**Table 3(a).** Confusion Matrix of Speaker IdentificationEngine (for MFCC features) in open set condition

Reference Speaker/ Test Speaker	51	55	S10	\$15	520	\$25	530	\$35	540	\$45	\$50	<u>555</u>	560	<u>\$65</u>	\$70	\$75	580	\$85	<u>\$90</u>	<u>\$95</u>
51	52.6	42.8	67.2	43.4	52.2	51.0	43.4	46.9	51.3	51.0	46.9	51.2	50.9	34.8	51.1	42.8	62.8	50.6	42.2	44.7
52	66.2	72.8	66.5	62.4	51.2	32.6	64.8	52.2	69.6	62.2	82.4	65.3	32.4	72.2	65.4	70.4	69.2	71.1	54.7	81.4
53	42.8	55.5	56.4	52.9	62.5	53.4	43.7	61.3	46.6	52.4	54.8	32.7	42.8	52.9	53.4	61.6	50.2	51.4	32.7	42.8
54	32.6	54.3	53.4	<mark>71.9</mark>	43.2	51.3	43.2	46.4	62.3	52.9	61.2	73.4	69.0	71.4	61.5	61.0	72.4	66.5	52.9	55.5
55	42.8	33.4	61.5	72.8	<mark>76.2</mark>	62.4	57.1	48.9	56.3	61.2	54.5	34.4	41.2	53.6	44.2	54.3	54.9	32.1	54.6	42.3
56	71.1	32.6	83.4	71.6	62.4	76.5	62.2	31.7	51.3	72.2	64.3	53.2	81.0	69.2	61.5	54.3	65.2	69.6	71.7	62.6
57	76.4	64.8	52.4	62.3	57.1	51.9	61.1	60.0	42.4	62.3	54.1	62.1	61.0	72.0	53.2	61.8	52.3	60.0	53.3	61.0
58	56.5	43.2	44.4	32.3	61.7	31.7	70.0	63.2	62.4	73.9	51.7	59.0	54.5	54.3	62.7	63.4	51.8	58.3	61.7	56.3
59	62.4	49.6	48.4	51.3	49.0	51.3	42.4	61.3	52.4	50.0	50.3	51.5	50.9	51.3	43.3	48.6	52.4	43.9	52.1	55.3
510	71.8	62.2	79.2	52.9	71.3	72.2	71.3	73.9	77.6	74.2	54.1	70.9	53.3	64.7	71.6	70.4	72.3	59.2	64.3	48.9
511	62.4	61.6	54.8	61.2	66.5	64.3	53.3	51.7	64.3	54.1	68.8	61.7	52.7	71.0	62.4	58.9	69.6	67.2	71.6	62.3
512	61.2	65.3	66.3	61.8	34.4	52.7	65.3	67.1	71.9	70.9	61.7	70.2	66.6	66.2	58.3	68.9	66.9	62.1	62.5	66.3
513	81.2	32.4	61.6	69.0	59.6	81.0	61.4	65.7	66.6	53.3	52.7	78.6	<u>69.2</u>	51.1	65.4	66.3	42.7	49.6	52.4	65.3
514	44.8	72.2	62.4	64.3	43.2	45.6	32.7	57.8	62.3	64.7	71.0	45.5	51.0	64.8	62.3	54.3	54.4	55.3	56.3	54.4
\$15	61.1	60.4	58.9	54.8	34.7	45.9	53.2	32.8	58.3	54.5	61.3	32.7	54.3	45.9	62.1	85.4	31.8	57.7	45.6	43.9
\$16	42.8	56.4	62.3	61.0	58.2	54.3	54.9	43.1	57.7	54.8	58.8	52.9	61.3	51.9	52.3	<u>59.7</u>	41.8	45.8	46.2	48.3
517	32.8	35.5	41.6	42.1	33.3	46.4	42.3	41.5	32.6	30.0	49.2	39.3	32.7	30.0	41.0	41.6	43.2	45.8	31.9	34.7
518	55.6	71.1	61.9	66.5	70.2	69.6	72.1	58.3	54.6	59.2	67.2	68.4	49.6	54.2	66.4	72.4	85.8	<u>79.0</u>	78.8	62.3
\$19	42.2	54.7	71.4	52.9	56.1	71.7	53.3	61.7	54.8	64.3	64.9	62.5	52.4	59.6	63.3	66.5	62.8	59.1	66.5	62.0
520	64.7	70.4	60.4	55.5	58.3	62.6	61.0	57.9	72.6	48.9	70.5	71.0	65.3	61.4	62.3	63.2	54.7	55.0	52.3	71.1

**Table 3(b).** Confusion matrix of speaker identification engine (for P features) in open set condition

Reference Speaker/ Test Speaker	51	55	\$10	\$15	520	\$25	\$30	\$35	540	\$45	S50	\$55	<i>\$60</i>	\$65	\$70	\$75	\$80	585	<u>\$90</u>	\$95
51	62.8	52.3	61.2	52.4	52.2	61.0	59.4	56.5	51.3	51.0	59.9	61.3	50.4	60.2	62.8	61.3	62.8	50.6	42.2	44.7
52	66.2	76.3	66.5	62.4	51.2	32.6	64.8	52.2	69.6	62.2	72.4	65.3	32.4	72.2	65.4	70.4	69.2	71.1	54.7	72.3
53	52.8	56.4	53.9	57.6	62.5	54.8	57.2	51.3	53.4	59.2	54.8	32.7	42.8	52.9	53.4	61.6	50.2	51.4	32.7	42.8
54	32.6	54.3	53.4	76.9	43.2	51.3	43.2	46.4	62.3	52.9	61.2	73.4	69.0	71.4	61.5	61.0	72.4	66.5	52.9	55.5
55	42.8	33.4	61.5	72.8	<mark>76.2</mark>	62.4	57.1	48.9	56.3	61.2	54.5	34.4	41.2	53.6	44.2	54.3	54.9	32.1	54.6	42.3
56	71.1	32.6	80.3	71.6	62.4	<u>76.5</u>	62.2	31.7	51.3	72.2	64.3	53.2	81.0	69.2	61.5	54.3	65.2	69.6	71.7	62.6
57	76.4	64.8	52.4	62.3	57.1	51.9	77.1	60.0	42.4	62.3	54.1	62.1	61.0	72.0	53.2	61.8	52.3	60.0	53.3	61.0
58	56.5	43.2	44.4	32.3	61.7	31.7	70.0	7 <u>3.2</u>	62.4	73.9	51.7	59.0	54.5	54.3	62.7	63.4	51.8	58.3	61.7	56.3
59	62.4	49.6	48.4	51.3	49.0	51.3	42.4	61.3	72.4	50.0	50.3	51.5	50.9	51.3	43.3	48.6	52.4	43.9	52.1	55.3
510	71.8	62.2	70.2	52.9	71.3	72.2	71.3	73.9	77.6	78.2	54.1	70.9	53.3	64.7	71.6	70.4	72.3	59.2	64.3	48.9
511	62.4	61.6	54.8	61.2	66.5	64.3	53.3	51.7	64.3	54.1	<mark>78.8</mark>	61.7	52.7	71.0	62.4	58.9	69.6	67.2	71.6	62.3
512	61.2	65.3	66.3	61.8	34.4	52.7	65.3	67.1	71.0	70.9	61.7	7 <u>1.2</u>	66.6	66.2	58.3	68.9	66.9	62.1	62.5	66.3
\$13	71.2	32.4	61.6	69.0	59.6	71.0	61.4	65.7	66.6	53.3	52.7	71.6	72.2	51.1	65.4	66.3	42.7	49.6	52.4	65.3
514	44.8	72.2	62.4	64.3	43.2	45.6	32.7	57.8	62.3	64.7	71.0	45.5	51.0	74.8	62.3	54.3	54.4	55.3	56.3	54.4
\$15	61.1	60.4	58.9	54.8	34.7	45.9	53.2	32.8	58.3	54.5	61.3	32.7	54.3	45.9	72.1	65.4	31.8	57.7	45.6	43.9
\$16	42.8	56.4	62.3	61.0	58.2	54.3	54.9	43.1	57.7	54.8	58.8	52.9	61.3	51.9	52.3	<u>69.7</u>	41.8	45.8	46.2	48.3
517	32.8	35.5	41.6	42.1	33.3	46.4	42.3	41.5	32.6	30.0	49.2	39.3	32.7	30.0	41.0	41.6	63.2	45.8	31.9	34.7
518	55.6	71.1	61.9	66.5	70.2	69.6	72.1	58.3	54.6	59.2	67.2	68.4	49.6	54.2	66.4	72.4	82.4	80.0	78.8	62.3
\$19	42.2	54.7	71.4	52.9	56.1	71.7	53.3	61.7	54.8	64.3	64.9	62.5	52.4	59.6	63.3	66.5	62.8	59.1	76.5	62.0
520	64.7	70.4	60.4	55.5	58.3	62.6	61.0	57.9	72.6	48.9	70.5	71.0	65.3	61.4	62.3	63.2	54.7	55.0	52.3	73.2

On the basis of confusion matrix, the speakers can be considered as nearest and farthest speakers (having highest and lowest accuracy).

Table 4.	Performance	of nearest	and	farthest	speakers
----------	-------------	------------	-----	----------	----------

Nearest	Accuracy	Farthest	Accuracy
Speakers	(%)	Speakers	(%)
5 and 50	92.2	5 and 60	42.2
55 and 60	78.9	20 and 55	34.8
80 and 85	86	-	-

Recognition performance has been also computed by comparing speech samples of 43 male and 57 female speakers as shown in Figure 6 (a, b).



**Figure 6(a).** Performance of male and female speakers based on MFCC features.



**Figure 6(b).** Performance of male and female speakers based on PLP features.

## 5. Conclusions

- Performance of PLP is better than the MFCC, particularly in the case of open set.
- The performance increases as the number of utterances in training samples increases.
- Performance of male speakers is found better than female speakers.
- On comparison with the earlier experimental results conducted by different researchers, it is

concluded that here the sentence level recognition give about

• 91.2% accuracy for MFCC and 92.3% for PLP which is better than the earlier experiments<sup>16</sup>.

## 6. References

- 1. Rabiner LR, Juang BH. Fundamentals of speech recognition. PTR Prentice-Hall; 1993.
- Agrawal SS, Bansal S, Pandey D, Tayal H. A Hidden Markov Model (HMM) based speaker identification system using mobile phone database of NATO (National Atlantic Treaty Organization) Words; ICA-Montreal, Canada. 2013.
- 3. Yuan J, Liberman M. Speaker identification on the Scotus Corpus. Available from: http://www.ling.upenn. edu/~jiahong/publications/c09.pdf
- 4. 4. Stern RM, Ponce P, Singh R. Feature extraction for robust automatic speech recognition using synchrony of zero crossings. Available from: http://citeseerx.ist.psu.edu/view-doc/download?doi=10.1.1.407.1507&r ep=rep1&type=pdf
- HTK Hidden Markov Toolkit Version 1.4 Manual. Cambridge University Engineering Department, Speech Group; Cambridge. 1992.
- Viterbi algorithm. Available from: http://en.wikipedia.org/ wiki/Viterbi\_algorithm
- Reddy V, Das PK. A HMM based text-prompted speaker verification system using HTK. Oriental COCOSDA; 2004 Nov 17-19. p. 313–6.
- Shrawankar U, Thakare V. Techniques for feature extraction in speech recognition system: A comparative study. Available from: http://arxiv.org/ftp/arxiv/papers/1305/ 1305.1145.pdf

- 9. Skowronski MD, Harris JG. Human factor cepstral coefficients. Cancun, Mexico; 2002 Dec.
- Aida–Zade KR, Ardil C, Rustamov SS. Investigation of combined use of MFCC and LPC features in speech recognition systems. Proceedings of World Academy of Science, Engineering and Technology; 2006 May; 13. ISSN 1307-6884.
- 11. Hermansky H. Perceptual Linear Predictive (PLP) analysis of speech. J Acous L Soc Am. 1990 Apr; 87(4).
- Hermansky H, Morgan N, Bayya A, Kohn P. RASTA-PLP speech analysis. Technical Report (TR-91-069). International Computer Science Institute; Berkeley, CA. 1991.
- Davis S, Mermelstein P. Comparison of parametric representations for monosyllabic word recognition in continuously spoken sentences. Published in Processing, IEEE Transactions on Acoustics, Speech and Signal. 1980 Aug; 28(4):357–66.
- Gauvain JL, Lamel LF, Adda G, Adda- Decker M. (1994): Speaker-independent continues speech dictation. Speech Communication. 1994; 15:21–37.
- Gupta R. Speech recognition for Hindi. M. Tech. Project Report. Department of Computer Science and Engineering, Mumbai, India: Indian Institute of Technology Bombay; 2006.
- 16. Kral P, Jezek K, Jedlicka P. Evaluation of automatic speaker recognition approaches. WESPAC X; Beijing, China. 2009.
- 17. Petr Motlicek. Feature extraction in speech coding and recognition. Report of PhD Research Internship in ASP Group, OGI-OHSU; 2001/2002
- Useful links. 1. http://www.nist.gov/speech, 2. http:// www. voxforge.ac.in, 3. http:// www.phon.ucl.ac.uk