

# Novel Quality Metric for Duration Variability Compensation in Speaker Verification using i-Vectors

Arnab Poddar\*, Md Sahidullah†, Goutam Saha§

\*§ Dept of Electronics and Electrical Communication Engineering, Indian Institute of Technology, Kharagpur, India

†Speech and Image Processing Unit, School of Computing, University of Eastern Finland, Joensuu, Finland

Email: \*arnabpoddar@iitkgp.ac.in, †sahid@cs.uef.fi, §gsaha@ece.iitkgp.ernet.in

**Abstract**—Automatic speaker verification (ASV) is the process to recognize persons using voice as biometric. The ASV systems show considerable recognition performance with sufficient amount of speech from matched condition. One of the crucial challenges of ASV technology is to improve recognition performance with speech segments of short duration. In short duration condition, the model parameters are not properly estimated due to inadequate speech information, and this results poor recognition accuracy even with the state-of-the-art i-vector based ASV system. We hypothesize that considering the estimation quality during recognition process would help to improve the ASV performance. This can be incorporated as a quality measure during fusion of ASV systems. This paper investigates a new quality measure for i-vector representation of speech utterances computed directly from Baum-Welch statistics. The proposed metric is subsequently used as quality measure during fusion of ASV systems. In experiments with the NIST SRE 2008 corpus, We have shown that inclusion of proposed quality metric exhibits considerable improvement in speaker verification performance. The results also indicate the potentiality of the proposed method in real-world scenario with short test utterances.

**Index Terms**—Short-segments, Duration Variability, Baum-Welch Statistics, Quality Measure, GMM-UBM, i-vector, Fusion, Speaker Recognition.

## I. INTRODUCTION

Automatic speaker verification (ASV) is a biometric recognition system where the voice is used as the trait [1], [2]. ASV is a convenient and non-invasive technology that can potentially be applied to various important applications, covering access of control, authentication of secure transactions over a telephone connection and forensic identification of suspects using voice samples [1], [3]. Contrasting to other biometrics, speaker recognition is a non-obtrusive technology and does not involve special purpose acquisition hardware other than a microphone. Even though speaker recognition research has been ongoing for more than four decades, the state-of-the-art speaker recognition systems still have several limitations[3], [4], [5].

Although state-of-the-art i-vector based ASV systems exhibit satisfactory performance with adequate speech data, but practically, the performance of such systems decline with limited duration data [6], [3], [4]. ASV system, in real-world applications requires satisfactory performance with short

duration speech which remains as an opportunity to explore further. The work in [7] attempted to model the duration variability in short duration as noise and also compensated with synthetically generated supporting i-vectors for speaker modeling. The work in [8] proposed to estimate the variability originated due to shorter utterances in i-vector space. The ASV systems suffers from the duration variability due to mismatch in train-test segments.

In the modern i-vector based ASV systems, Baum-Welch (BW) statistics are indispensable intermediate parameters, which totally represent the extracted speech-features. The quality of estimation of BW statistics is degraded in short duration, which introduces sparsity due to insufficient data. The sparsity arises due to shortage of speech data as it fails to update most of the Gaussian components in BW statistics [9], [10]. The present de-facto ASV systems do not include the information regarding the quality of speaker model estimation. We consider BW statistics not only as the intermediate parameters for speaker model estimation, but also as a source to determine quality of speaker model estimation.

In this work, we introduce a metric to measure the quality of intermediate ASV system parameters. This work contributes to incorporate the information regarding quality of the speaker model estimation for the first time in ASV to the best of our knowledge. the proposed metric is estimated directly from the intermediate system parameters of i-vector based ASV system. This metric attempt to represent the impact of duration on intermediate BW statistics by calculating the difference between intermediate BW statistics and universal background parameters. The proposed dissimilarity metric do not require additional parameters to be estimated and require negligible computation cost as intermediate statistics are inherently calculated by state-of-the-art ASV systems. In the classification module of ASV systems, *Gaussian mixture model-universal background model* (GMM-UBM) [2] and i-vector [11] are used widely. An exhaustive comparison of the two techniques, including the short duration effect, reveal that though i-vector outperforms the GMM-UBM for longer speech utterances, but the GMM-UBM is considerably relevant for short duration condition [4]. The observation inspire us to fuse classifiers.

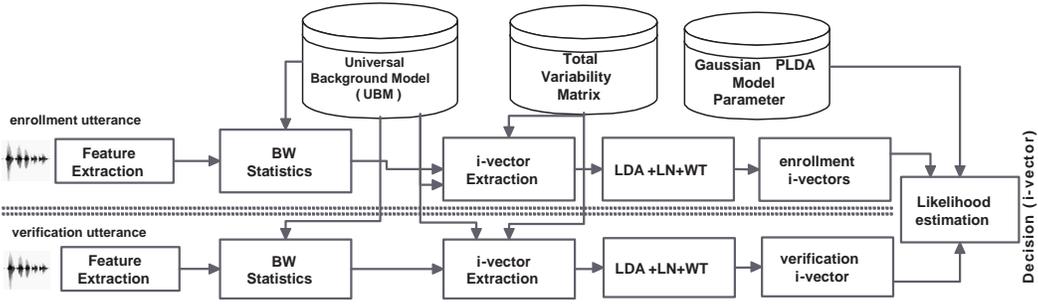


Fig. 1. Block diagram for i-vector based ASV system.

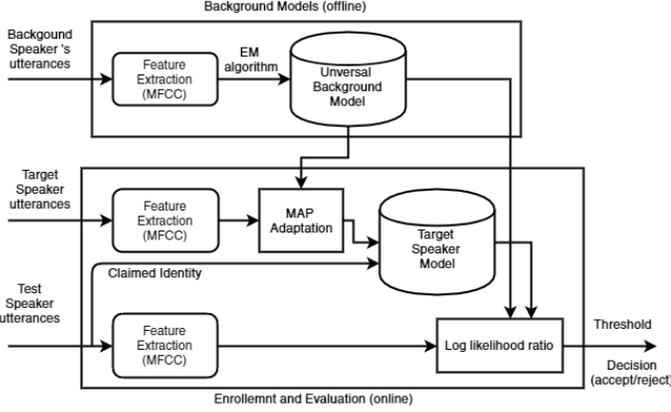


Fig. 2. Block diagram for GMM-UBM based ASV system.

Additionally, the proposed similarity metric is incorporated in fusion stage as quality information of speech to compensate the short duration effect. Incorporation of quality measures not only showed considerable improvement in performance in various duration conditions. The proposed systems showed more improvement for practical requirement i.e., in short duration cases.

In the rest of the paper, theoretical aspects of widely used i-vector and GMM-UBM based ASV techniques are illustrated in Section II. An analysis on intermediate parameters is presented in section III. Subsequently, Section IV and V discuss the proposed quality aided fusion based system and experimental results. Finally, we conclude the paper in Section VI.

## II. ASV SYSTEMS

Here, we discuss two popularly implemented ASV techniques, namely GMM-UBM [2] and i-vector representation of speech utterance [11]. Fig. 1 shows i-vector based the ASV system and Fig. 2 represents the GMM-UBM system.

### A. GMM-UBM system

In GMM-UBM approach, initially, a GMM is estimated with a large volume of voice from a large number of speakers who may not participate in the verification process [2]. The estimated background model is termed as *universal background model* (UBM), written as  $\lambda_{\text{UBM}} = \{w_i, \bar{\mu}_i, \bar{\Sigma}_i; i = 1, 2, \dots, K\}$ . Here,  $K$  represents the number of Gaussian

components in the mixture,  $w_i$  stands for the prior weights of the  $i$ -th Gaussian mixture components,  $\mu_i$  represents the mean and  $\Sigma_i$  represents the co-variance matrix. The parameter  $w_i$  meets the condition  $\sum_{i=1}^K w_i = 1$ .

The  $S$  speakers' GMM models are mathematically represented as  $\{\lambda_1, \lambda_2, \dots, \lambda_S\}$ . We estimate the enrollment speaker model by adapting of the UBM model parameters using *maximum-a-posteriori* (MAP) technique [2]. Initially, sufficient statistics  $N_i$  (zero order),  $\mathbf{E}_i$  (1st order) and  $\mathbf{F}_i$  (2nd order) from a enrollment speaker's utterance with  $C$  active frames  $\mathbf{X} = \{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_C\}$ , are computed as,

$$N_i = \sum_{t=1}^C Pr(i|\mathbf{x}_t), \quad (1)$$

$$\mathbf{E}_i(\mathbf{X}) = \frac{1}{N_i} \sum_{t=1}^C Pr(i|\mathbf{x}_t)\mathbf{x}_t, \quad (2)$$

Here the distribution of probability of Gaussian mixture components  $Pr(i|\mathbf{x}_t)$  for given speech segments with  $C$  frames  $\mathbf{X}^{\text{train}} = \{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_C\}$  is formulated by

$$Pr(i|\mathbf{x}_t) = \frac{w_i p_i(\mathbf{x}_t)}{\sum_{j=1}^K w_j p_j(\mathbf{x}_t)} \quad (3)$$

where all probability density is a  $K$ -dimensional Gaussian variable of the form

$$p_i(\mathbf{x}) = \frac{1}{(2\pi)^{K/2} |\bar{\Sigma}_i|^{1/2}} \exp\left\{-\frac{1}{2}(\mathbf{x} - \bar{\mu}_i)^\top \bar{\Sigma}_i^{-1}(\mathbf{x} - \bar{\mu}_i)\right\}. \quad (4)$$

Conventionally, only the mean parameters of speaker's GMM model are adapted to estimate the enrollment models. It makes the speaker model estimation process computationally efficient [2].

During evaluation, the log-likelihood ratio of verification feature vectors are computed,  $\mathbf{X}^{\text{test}} = \{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_C\}$  against both target-speaker model and the background model as,

$$\Lambda_{\text{GMM-UBM}}(\mathbf{X}^{\text{test}}) = \log p(\mathbf{X}^{\text{test}}|\lambda_{\text{target}}) - \log p(\mathbf{X}^{\text{test}}|\lambda_{\text{UBM}}) \quad (5)$$

Finally, a threshold ( $\theta$ ) is adjusted to determine whether the claimed identity will be *accepted* or *rejected*. If  $\Lambda_{\text{GMM-UBM}}(\mathbf{X}) \geq \theta$ , the claim is accepted, otherwise rejected.

## B. i-vector Extraction

The i-vectors transform the *GMM supervector* into a lower dimensional subspace [11]. The adapted GMM supervector of  $i$ -th speaker,  $\mathbf{m}_i$ , is transformed as,

$$\mathbf{m}_i = \bar{\mathbf{m}} + \Phi \mathbf{y}, \quad (6)$$

here  $\Phi$  is a matrix of lower-rank, denoting the channel and speaker independent subspace,  $\mathbf{y}$  is i-vector,  $\bar{\mathbf{m}}$  represents the channel and speaker independent supervector ( $\bar{\mathbf{m}}$ ). Initially,  $\Phi$  is estimated with large volume of voice utterances collected from various persons [11]. Subsequently the corresponding i-vectors are calculated with the 1<sup>st</sup> order and zeroth order BW statistics  $\mathbf{E}_i$  and  $N_i$ , respectively.

Initially, sufficient statistics  $N_i$  (zero order),  $\mathbf{E}_i$  (1<sup>st</sup> order) and  $\mathbf{F}_i$  (2<sup>nd</sup> order) from a speaker's voice segment, consisting of  $C$  active frames  $\mathbf{X} = \{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_C\}$ , are computed as,

$$N_i = \sum_{t=1}^C Pr(i|\mathbf{x}_t), \quad (7)$$

$$\mathbf{E}_i(\mathbf{X}) = \frac{1}{N_i} \sum_{t=1}^C Pr(i|\mathbf{x}_t) \mathbf{x}_t, \quad (8)$$

here the distribution of probability of Gaussian components  $Pr(i|\mathbf{x}_t)$  conditioned on given voice segment with  $C$  speech frames  $\mathbf{X}^{\text{train}} = \{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_C\}$  is given by

$$Pr(i|\mathbf{x}_t) = \frac{w_i p_i(\mathbf{x}_t)}{\sum_{j=1}^K w_j p_j(\mathbf{x}_t)} \quad (9)$$

where each component density is a  $d$ -variate Gaussian function of the form as shown in Eq. 4. We consider the prior distribution of i-vectors  $p(\mathbf{y})$  is normally distributed as  $\mathcal{N}(0, \mathbf{I})$ . The corresponding posterior distribution of  $p(\mathbf{E}|\mathbf{y})$ , is assumed as  $p(\mathbf{E}|\mathbf{y}) = \mathcal{N}(\Phi \mathbf{y}, \mathbf{N}^{-1} \Sigma)$ . The intermediate parameter  $\mathbf{N}$  is computed as a diagonal matrix having  $N_i$  as its diagonal elements [11]. The MAP estimate of  $(\mathbf{y}|\mathbf{E})$  is computed as

$$\mathbb{E}(\mathbf{y}|\mathbf{E}) = (\mathbf{I} + \Phi^\top \Sigma^{-1} \mathbf{N} \Phi)^{-1} \Phi^\top \Sigma^{-1} \mathbf{N} (\mathbf{E} - \bar{\mathbf{m}}) \quad (10)$$

The expectation of  $(\mathbf{y}|\mathbf{E})$  is termed as the i-vector of a given voice segment  $\mathbf{X}$  [11].

The verification scores in i-vector GPLDA framework, is calculated as the likelihood ratio [12]. For a verification trial, the projected verification and enrollment i-vectors  $\mathbf{z}_{\text{test}}$  and  $\mathbf{z}_{\text{target}}$  respectively are used to estimate the likelihood ratio  $\Lambda_{\text{GPLDA}}(\mathbf{z}_{\text{target}}, \mathbf{z}_{\text{test}})$  as,

$$\Lambda_{\text{GPLDA}}(\mathbf{z}_{\text{target}}, \mathbf{z}_{\text{test}}) = \log \frac{p(\mathbf{z}_{\text{target}}, \mathbf{z}_{\text{test}} | H_1)}{p(\mathbf{z}_{\text{target}} | H_0) p(\mathbf{z}_{\text{test}} | H_0)} \quad (11)$$

here  $H_1$  hypothesizes the projected i-vectors belong to the same person. On the other hand,  $H_0$  denotes the hypothesis where the i-vectors belong to different person.

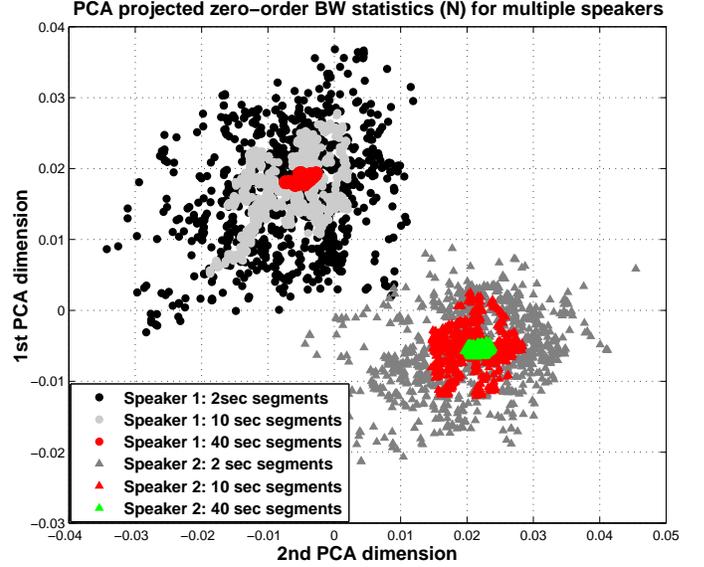


Fig. 3. Scatter plot of PCA projected NBS ( $\tilde{N}$ ) for two speakers.

## III. ANALYSIS OF AND CHARACTERISTICS OF BW STATISTICS

BW statistics represent the overall extracted information from the speech and are transformed into i-vectors using the pre-estimated *universal background model* (UBM) [11], [2].

Since BW statistics is an indispensable intermediate step in ASV, we investigate its characteristics in short duration. The zeroth order BW statistics ( $N_i$ ) is estimated as,  $N_i = \sum_{t=1}^C Pr(i|\mathbf{x}_t)$ , where, a speech segment with  $C$  frames is represented as  $\mathbf{X} = \{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_K\}$ , and  $Pr()$  is the prior probability of  $i$ -th Gaussian component. Summing over all Gaussian mixture components  $K$  we obtain,

$$\sum_{i=1}^K N_i = \sum_{i=1}^K \sum_{t=1}^C Pr(i|\mathbf{x}_t) = C. \quad (12)$$

The equation indicates that  $N_i$  is dependent on segment duration, i.e.,  $C$ . Normalizing  $N_i$  with number of frames ( $C$ ) we get

$$\tilde{N}_i = \frac{1}{C} \sum_{t=1}^C Pr(i|\mathbf{x}_t) \quad (13)$$

and  $\sum_{i=1}^K \tilde{N}_i = 1$  hence it has the same property as weights of the GMM UBM i.e.,  $\sum_{i=1}^K w_i = 1$ .

Now,  $\tilde{N}_i$  can be regarded as the mixture weights of the Gaussian components  $i$  for a particular speech segment  $s$ . It is a standard statistical hypothesis that intermediate BW statistics can be estimated more efficiently with sufficient volume speech corpus, which is likely to include all possible kinds of variability proportionately. Hence, we expect that large number of speech frames ( $C$ ) in the speech segment would be advantageous for improvement of quality of  $\tilde{N}_i$ .

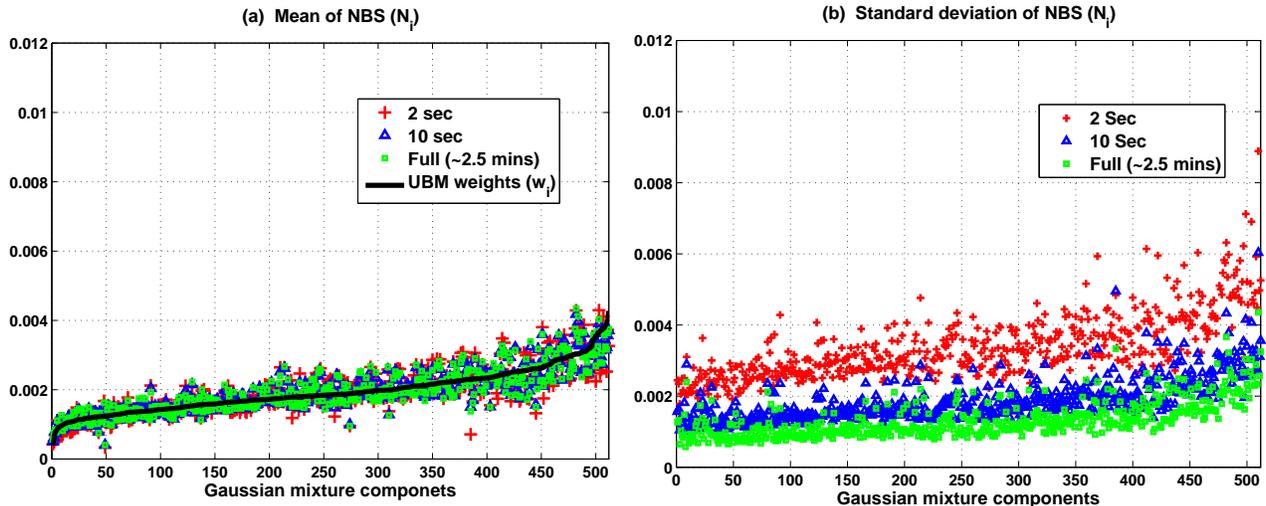


Fig. 4. Mean and standard deviation of NBS, calculated from voice utterances of 1270 male speakers (NIST 2008), are presented in (a) and (b) respectively. Means and standard deviation of NBS are plotted for three duration conditions, e.g., 2 Sec, 10 Sec and full length ( 2.5 min). Weights of GMM-UBM for corresponding Gaussian mixture component are shown in (a).

However, the intermediate statistics may be expected to be updated more sparsely due to reduced speech data or degraded quality of speech. On this core note, the characteristics of  $\tilde{N}_i$  are investigated.

Here we present analysis on NBS and its characteristics for different duration conditions, to observe the impact of duration. We considered the recordings from telephonic conversations in NIST SRE 2008 (*short2*) corpus. Truncated voice segments (40, 10, 2 sec) are considered for the analytical experiments. For truncation of the long segments, the initial speech frame is selected randomly and required number of successive active frames are pruned. Similarly, 1000 truncated utterances are generated for the duration conditions, under consideration.

We apply the principal component analysis (PCA) on the feature matrix. Subsequently, we show the major two projected components. In Fig.3, the projected components of different truncated segments of 2 Sec, 10 Sec and 40 Sec are shown. We estimate the matrix for PCA projection from the generated 1000 truncated segments. We observe that the NBS show greater variability in limited duration. Higher variability in NBS for limited duration, deteriorates the quality in i-vector model. The change in variability of BW statistics with duration of the speech segments indicate that BW statistics is associated with duration or estimation quality.

Fig. 4 (a) and (b) represents the mean and standard deviation of NBS ( $\tilde{N}_i$ ) respectively, for three duration conditions (2 sec, 10 sec and full length). The means of NBS ( $\tilde{N}_i$ ) are calculated using 1270 male speakers from NIST 2008 telephone corpus. The weights of GMM-UBM of corresponding mixture components ( $w_i$ ) are presented simultaneously in Fig. 4 (a). The short segments in Fig. 4 (b) showed greater standard deviation referring greater variability introduced in NBS ( $\tilde{N}_i$ ). We observe gradual increment in variability when the length of speech segments are shortened. As the variability in NBS is affected

by duration, we hypothesize that the information of duration variability can be estimated from NBS and also can be treated as the source of the information about the quality of speaker-model estimation. It is observed in Fig. 4 (a) that the means of  $\tilde{N}$  for different duration condition follows the value of UBM weight of corresponding Gaussian mixture component ( $w_i$ ). We also observe that the means of different duration condition for a particular Gaussian component remains nearly equal. This is observed in almost all Gaussian components shown in Fig. 4 (a). These observations on means of  $\tilde{N}$  distribution and weights of corresponding Gaussian mixture component inspired us to use GMM-UBM weights ( $w_i$ ) as reference to measure the variability in  $\tilde{N}$ .

#### IV. PROPOSED QUALITY MEASURE AND ASV SYSTEM FUSION

The observations in previous section demonstrate that the sparsity in BW statistics is associated with duration of speech and quality of speaker model estimation. The sparsity in  $\tilde{N}$  increases in short duration, indicating lower quality of estimation. From the observations in Section III, we consider that the BW statistics not only as intermediate parameters but also the source of estimating the quality of speaker model estimation. Here, we propose to quantify the dissimilarity between normalized zeroth order BW statistics  $\tilde{N}_i$  and prior of corresponding Gaussian component of UBM model  $w_i$ . We further use it as a quality metric. Subsequently, it is incorporated as supporting information in proposed ASV system. The mathematical expressions to model the quality  $Q$  of a segment  $s$  is given by

$$Q_s(\tilde{N}_s) = \sum_{i=1}^K |\tilde{N}_{i,s} - w_{i,ubm}| \quad (14)$$

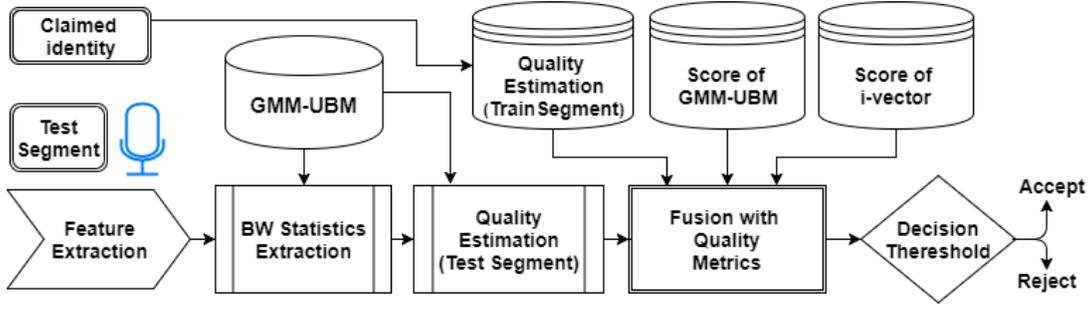


Fig. 5. Block diagram for the proposed quality estimation and quality metric incorporated fusion based approach for ASV.

TABLE I  
DETAILS OF SPEECH CORPUS AND CEPSTRAL FEATURES.

Specifications	#target model	#test segments	#genuine trials	#imposter trials
<b>NIST 2008</b>	442	854	874	11637
<b>specifications: Features and Development parameters</b>				
<b>MFCC</b>	Dimension: 19+19 $\Delta$ +19 $\Delta\Delta$ ; 20ms Hamming			
<b>GMM-UBM</b>	Dimension: 512 Data: NIST SRE '04, '05, Switchboard II			
<b>TV (<math>\Phi</math>) Matrix</b>	Dimension: 400; Data: NIST SRE '04, '05, '06, Switchboard II			
<b>GPLDA</b>	Dimension: 150; Data: NIST SRE '04, '05, '06, Switchboard II			

The proposed quality metric attempts to measure dissimilarity of  $N_i$  from the weights of UBM ( $w_i$ ) which is treated as reference.

ASV systems based on fusion approach have found wide applications [13], [14]. Though i-vector [11], [8] and GMM-UBM [2] based ASV systems have different modeling approaches, they exhibit similar performance in short utterance cases [4]. Here, we exploit the information captured simultaneously by GMM-UBM and modern i-vector GPLDA using linear fusion and subsequently incorporate the quality metric  $Q$  as additional information. The fusion parameters are trained using logistic regression objective using the BOSARIS toolkit [15]. We confine our work to score level fusion with fusion function  $f$  which combines two base classifier score  $\Lambda_{\text{UBM}}$  and  $\Lambda_{\text{i-vector}}$  into a single match score  $\mathbf{\Lambda} = \{\Lambda_{\text{UBM}}, \Lambda_{\text{i-vector}}\}^T$ . The decision is made by a predefined score threshold  $\theta$ . The trained linear fusion classifier is of the form

$$f_{\alpha,\theta}(\mathbf{\Lambda}) = \alpha^T \mathbf{\Lambda} + \theta \quad (15)$$

The fusion function, incorporating quality measure  $Q$  is represented by:

$$f_Q(\mathbf{\Lambda}) = \alpha^T \mathbf{\Lambda} + \theta + \beta * Q(\tilde{N}_{\text{enrollment}})Q(\tilde{N}_{\text{verification}}) \quad (16)$$

A speaker is accepted if and only if  $f_{Q,\alpha,\beta,\theta}(\mathbf{\Lambda}) \geq 0$ . When for a quality fusion classifier  $f_Q(\mathbf{\Lambda})$  with parameters  $(\alpha, \beta, \theta)$ , the development data  $D$  and an empirical cost function  $\hat{C}((\alpha, \beta, \theta), D)$  are given, the optimal fusion device is ob-

TABLE II  
RESULTS OF FUSION OF GMM-UBM AND I-VECTOR BASED SYSTEM WITH PROPOSED QUALITY METRIC ON NIST 2008 *Truncated Train - Truncated Test* TELEPHONE CORPORA

Train-Test duration	Metric	GMM UBM	i-vect GPLDA	linear fusion	Quality fusion
2s-2s	EER	35.24	36.84	33.05	<b>31.92</b>
	DCF	9.69	9.93	9.54	<b>9.50</b>
5s-5s	EER	25.25	24.37	23.11	<b>21.25</b>
	DCF	8.89	8.65	8.27	<b>8.13</b>
10s-10s	EER	14.98	14.53	14.05	<b>13.15</b>
	DCF	6.58	6.40	6.25	<b>5.86</b>

TABLE III  
RESULTS OF FUSION OF GMM-UBM AND I-VECTOR BASED SYSTEM WITH PROPOSED QUALITY METRIC ON NIST 2008 *Long Train- Truncated Test* TELEPHONE CORPORA

Train-Test duration	Metric	GMM UBM	i-vect GPLDA	linear fusion	Quality fusion
Full-2s	EER	21.56	19.67	19.10	<b>16.81</b>
	DCF	7.75	7.91	7.53	<b>7.01</b>
Full-5s	EER	17.73	13.50	12.23	<b>11.67</b>
	DCF	7.32	5.99	5.69	<b>5.42</b>
Full-10s	EER	16.66	9.29	9.72	<b>9.09</b>
	DCF	6.75	4.50	4.59	<b>4.28</b>

tained by  $(\alpha^{dev}, \beta^{dev}, \theta^{dev}) = \text{argmin}_{\alpha,\beta,\theta} \hat{C}((\alpha, \beta, \theta), D)$ . Here, the *decision cost function* is adopted as

$$C_{\text{det}}(\theta) = C_{\text{miss}}P_{\text{miss}}(\theta)P_{\text{tar}} + C_{\text{fa}}P_{\text{fa}}(\theta)(1 - P_{\text{tar}}) \quad (17)$$

where  $P_{\text{tar}}$  is the prior probability of an original speaker,  $C_{\text{miss}}$  is the cost of a miss and  $C_{\text{fa}}$  is the cost of false alarm. A diagrammatic representation of the proposed fusion based approach to include the quality metrics is presented in Fig. 5.

## V. EXPERIMENTAL RESULTS AND DISCUSSION

Here, Mel-frequency cepstral coefficients (MFCC) appending delta ( $\Delta$ ) and double-delta ( $\Delta\Delta$ ) coefficients are used for experiments [16]. The non-speech frames are discarded by a spectrum energy-based detector (SAD) and finally, cepstral mean and variance normalization (CMVN) is applied as feature normalization [17], [18]. A gender-specific UBM (male) is used. We conducted the experiments on NIST speaker recognition evaluation (SRE) corpus 2008. We consid-

ered *short2-short3* task<sup>1</sup> on *telephone-telephone* part of male speakers. The details of NIST SE 2008 are given in Table I. Channel compensation for i-vectors are done using Gaussian probabilistic linear discriminant analysis (GPLDA) [12], [6], [4]. A brief synopsis of the development parameters used in the experiments are outlined in Table I. To generate short utterances, we truncate the long speech utterances in 2 sec (200 active frames), 5 sec (500 active frames), 10 sec (1000 active frames) duration removing prior 500 active speech frames to avoid phonetic similarity in initial greetings to avoid text-dependence [9], [19].

Quality measures of speech signals are used to support the fusion based ASV system. Performance measures of GMM-UBM, i-vector and fusion, using quality measures are depicted separately in Table II and Table III. Separate experiments are conducted with long enrollment data (III) and also short enrollment data (II). A total of six different duration conditions are used for experiments as shown in Table II and III. The results are shown in both *equal error rate* (EER) and *minimum detection cost function* (minDCF) [6]. Incorporation of quality metric exhibited considerably high relative improvement over state-of-the-art, in conditions like *full-2 sec*, *full-5 sec*, *5 sec-5 sec*, *2 sec-2 sec* etc. These conditions are more close to desirable real-time requirements of ASV systems which encourages to find implementations of proposed system. Consistent improvement of accuracy of the ASV system in various duration established relevance of the proposed quality measures based on intermediate statistics. The system is more suitable when the duration of speech utterances are limited, especially when it is trained with long enrollment data and tested with very short duration of speech.

## VI. CONCLUSION

This work investigates a new metric for measuring the quality of the i-vector estimation process. The metric is formulated using the Baum-Welch statistics and UBM parameters. The proposed metric helps to improve the ASV performance when incorporated as a side information during system combinations. The relative improvement is considerably more when tested with short test data. This quality metric requires no additional parameters to be estimated. In our current work, we have proposed a simple scheme where the absolute differences of BW statistics and UBM parameters are used for measuring the quality. Further investigation can be conducted by adopting other dissimilarity metrics with a goal to find the optimum one. Other possible future directions include evaluating the performance in noisy conditions, compatibility test of the proposed approach with other features such as deep neural network based bottleneck features, etc.

## ACKNOWLEDGMENT

The authors take the opportunity to acknowledge Indian Space Research Organization (ISRO) for financing the research partially. The authors also express gratitude to Mr.

Monisankha Pal and Mrs. Shefali Waldekar for technical discussions and grammatical corrections respectively.

## REFERENCES

- [1] T. Kinnunen and H. Li, "An overview of text-independent speaker recognition: From features to supervectors," *Speech communication*, vol. 52, no. 1, pp. 12–40, 2010.
- [2] D. A. Reynolds, T. F. Quatieri, and R. B. Dunn, "Speaker verification using adapted Gaussian mixture models," *Digital signal processing*, vol. 10, no. 1, pp. 19–41, 2000.
- [3] A. Poddar, M. Sahidullah, and G. Saha, "Speaker verification with short utterances: A review of challenges, trends and opportunities," *IET Biometrics*, 2017.
- [4] —, "Performance comparison of speaker recognition systems in presence of duration variability," in *2015 Annual IEEE India Conference (INDICON)*. IEEE, 2015, pp. 1–6.
- [5] L. Li, D. Wang, C. Zhang, and T. F. Zheng, "Improving short utterance speaker recognition by modeling speech unit classes," *IEEE/ACM Transactions on Audio, Speech, and Language Processing*, vol. 24, no. 6, pp. 1129–1139, 2016.
- [6] A. Kanagasundaram, R. Vogt, D. B. Dean, S. Sridharan, and M. W. Mason, "I-vector based speaker recognition on short utterances," in *Proceedings of INTERSPEECH*. International Speech Communication Association (ISCA), 2011, pp. 2341–2344.
- [7] T. Hasan, R. Saeidi, J. H. Hansen, D. van Leeuwen *et al.*, "Duration mismatch compensation for i-vector based speaker recognition systems," in *Acoustics, Speech and Signal Processing (ICASSP), 2013 IEEE International Conference on*. IEEE, 2013, pp. 7663–7667.
- [8] A. Kanagasundaram, D. Dean, S. Sridharan, J. Gonzalez-Dominguez, J. Gonzalez-Rodriguez, and D. Ramos, "Improving short utterance i-vector speaker verification using utterance variance modelling and compensation techniques," *Speech Communication*, vol. 59, pp. 69–82, 2014.
- [9] A. Poddar, M. Sahidullah, and G. Saha, "An adaptive i-vector extraction for speaker verification with short utterance," in *Proc. International Conference on Pattern Recognition and Machine Intelligence (PRMI)*. Springer, 2017.
- [10] W. Li, T. Fu, H. You, J. Zhu, and N. Chen, "Feature sparsity analysis for i-vector based speaker verification," *Speech Communication*, vol. 80, pp. 60–70, 2016.
- [11] N. Dehak, P. Kenny, R. Dehak, P. Dumouchel, and P. Ouellet, "Front-end factor analysis for speaker verification," *Audio, Speech, and Language Processing, IEEE Transactions on*, vol. 19, no. 4, pp. 788–798, 2011.
- [12] P. Kenny, "Bayesian speaker verification with heavy-tailed priors," in *Odyssey*, 2010, p. 14.
- [13] T. Hasan, S. O. Sadjadi, G. Liu, N. Shokouhi, H. Boril, and J. H. Hansen, "CRSS systems for 2012 NIST speaker recognition evaluation," in *ICASSP*. IEEE, 2013, pp. 6783–6787.
- [14] V. Hautamaki, T. Kinnunen, F. Sedláč, K. A. Lee, B. Ma, and H. Li, "Sparse classifier fusion for speaker verification," *Audio, Speech, and Language Processing, IEEE Transactions on*, vol. 21, no. 8, pp. 1622–1631, 2013.
- [15] N. Brümmer and E. De Villiers, "The bosaris toolkit: Theory, algorithms and code for surviving the new dcf," *arXiv preprint arXiv:1304.2865*, 2013.
- [16] M. Sahidullah and T. Kinnunen, "Local spectral variability features for speaker verification," *Digital Signal Processing*, vol. 50, pp. 1–11, 2016.
- [17] M. Sahidullah and G. Saha, "Design, analysis and experimental evaluation of block based transformation in MFCC computation for speaker recognition," *Speech Communication*, vol. 54, no. 4, pp. 543–565, 2012.
- [18] —, "A novel windowing technique for efficient computation of MFCC for speaker recognition," *Signal Processing Letters*, vol. 20, no. 2, pp. 149–152, 2013.
- [19] A. Poddar, M. Sahidullah, and G. Saha, "Improved i-vector extraction technique for speaker verification with short utterances," in *International Journal of Speech Technology*. Springer, 2017.

<sup>1</sup>[http://www.itl.nist.gov/iad/mig/tests/sre/2008/sre08\\_evalplan\\_release4.pdf](http://www.itl.nist.gov/iad/mig/tests/sre/2008/sre08_evalplan_release4.pdf)