

# The HKCUPU System for The NIST 2010 Speaker Recognition Evaluation

Weiwu JIANG<sup>1</sup>, Man-Wai MAK<sup>2</sup>, Wei RAO<sup>2</sup> and Helen MENG<sup>1</sup>

<sup>1</sup>Dept. of System Engineering and Engineering Management, The Chinese University of Hong Kong

<sup>2</sup>Dept. of Electronic and Information Engineering, The Hong Kong Polytechnic University

## Introduction

This paper presents the HKCUPU speaker recognition system submitted to NIST 2010 speaker recognition evaluation (SRE). The system comprises five subsystems, each with different acoustic features, session-variability reduction methods, speaker modeling and scoring methods and classifiers. This paper reports the results of individual and fusion systems for the core test and highlights the improvements made by our newly proposed JFA-Fishvoice (FSH) subsystem. Results show that FSH outperforms JFA when its projection matrix is channel-dependent (telephone or microphone) and that FSH is complementary to other state-of-the-art techniques. It was also found that VAD is an important pre-processing step for interview speech.

## Subsystems

Sub-System	Description
JFA	Joint Factor Analysis
JSV	Linear SVM on JFA-GMM supervectors
JSF	Cosine-kernel SVM on JFA speaker factors
GSV	Linear SVM on GMM-supervectors with NAP
FSH	Cosine distance on Fishvoice-projected speaker factors

## Voice Activity Detection

Type of Speech	VAD Method
Telephone	<ul style="list-style-type: none"> <li>Energy-based VAD for GSV subsystem</li> <li>ETSI Adaptive Multi-Rate (AMR) VAD for other four subsystems</li> </ul>
Microphone and Interview	<ul style="list-style-type: none"> <li>Spectral subtraction followed by energy-based VAD with crosstalk removal</li> </ul>

## Acoustic Features

Subsystem	Features and Dimension	Frame Size
JFA	17 MFCC <sub>0+Δ+ΔΔ</sub> (51Dim)	25ms
JSV	17 MFCC <sub>0+Δ+ΔΔ</sub> (51Dim)	25ms
JSF	12 PLP <sub>E+Δ+ΔΔ+ΔΔΔ</sub> (52 Dim)	20ms
GSV	12 MFCC <sub>Δ</sub> (24 Dim)	25ms
FSH	12 PLP <sub>E+Δ+ΔΔ+ΔΔΔ</sub> (52 Dim)	20ms

- GSV: Cepstral mean normalization followed by feature warping.
- Others: Mean-variance-normalization followed by feature warping.

## Subsystem Description

### JFA Subsystem:

- 2048-Gaussian GMM-supervector of speaker  $s$ :

$$\mathbf{m}(s) = \mathbf{m} + \mathbf{V}\mathbf{y}(s) + \mathbf{U}\mathbf{x}(s, h_s)$$

where  $\mathbf{m}$  is the UBM supervector,  $\mathbf{U}$  is the Eigenchannel matrix,  $\mathbf{V}$  is the Eigenvoice matrix,  $\mathbf{x}(s, h)$  is the speaker-dependent Eigenchannel factor,  $\mathbf{y}(s)$  is the session- and speaker-dependent Eigenvoice factor, and  $h_s$  represents the enrollment channel.

### JSV Subsystem:

- Use GMM-supervectors of target speaker  $s$  and claimant  $c$  obtained from JFA as feature vectors:

$$\mathbf{m}(s) = \mathbf{m} + \mathbf{V}\mathbf{y}(s)$$

$$\mathbf{m}(c) = \mathbf{m} + \mathbf{V}\mathbf{y}(c)$$

- Linear SVM Scoring:

$$S_{\text{SVM}}(s, c) = \alpha_{s,0} \langle \mathbf{m}(s), \mathbf{m}(c) \rangle - \sum_{i \in \text{SV}^-} \alpha_{s,i} \langle \mathbf{m}(b_i), \mathbf{m}(c) \rangle + d_s$$

where  $b_i$  is the  $i$ -th background speaker and  $\text{SV}^-$  contains the support vector indexes of background speakers.

### JSF Subsystem:

- Use speaker-factors of target speaker  $s$  and claimant  $c$  obtained from JFA as feature vectors:

$$\mathbf{m}(s) = \mathbf{m} + \mathbf{V}\mathbf{y}(s) + \mathbf{U}\mathbf{x}(s, h_s)$$

$$\mathbf{m}(c) = \mathbf{m} + \mathbf{V}\mathbf{y}(c) + \mathbf{U}\mathbf{x}(c, h_c)$$

- Cosine-kernel SVM Scoring:

$$S_{\text{CSVM}}(s, c) = \alpha_{s,0} K(\mathbf{y}(s), \mathbf{y}(c)) - \sum_{i \in \text{SV}^-} \alpha_{s,i} K(\mathbf{y}(b_i), \mathbf{y}(c)) + d_s$$

### GSV Subsystem:

- Use MAP adapted 512-Gaussian GMM-supervectors
- Transformed by NAP matrices (Corank = 16 for tel; Corank = 128 for mic/interview)
- Linear SVM Scoring with T-norm

### FSH Subsystem:

- Apply JFA speaker factors  $\mathbf{y}(s)$  as feature vectors to estimate a nonparametric Fisher discriminant projection matrix  $\mathbf{W}$ :

$$\mathbf{W} = \mathbf{W}_1 \mathbf{W}_2 \mathbf{W}_3$$

where  $\mathbf{W}_1$  is PCA projection matrix,  $\mathbf{W}_2$  is whitened within-class projection matrix, and  $\mathbf{W}_3$  is nonparametric between-class projection matrix.  $\mathbf{W}_3$  focuses on the boundaries between speakers without using parametric models to approximate the distribution of  $\mathbf{y}(s)$ .

$$\mathbf{W}_1 = \arg \max_{\mathbf{W}} \|\mathbf{W}^T \Sigma \mathbf{W}\| \quad \mathbf{W}_3 = \arg \max_{\mathbf{W}} \|\mathbf{W}^T \mathbf{S}_y \mathbf{W}\|$$

$$\mathbf{W}_2^T \mathbf{S}_w \mathbf{W}_2 = \mathbf{I} \Rightarrow \mathbf{W}_2 = \Phi \Lambda^{-\frac{1}{2}}$$

- Cosine distance scoring:

$$S_{\text{FSH}}(s, c) = \frac{\langle \mathbf{W}\mathbf{y}(s), \mathbf{W}\mathbf{y}(c) \rangle}{\|\mathbf{W}\mathbf{y}(s)\| \|\mathbf{W}\mathbf{y}(c)\|}$$

Fusion System	CC5	CC6	CC8
JFA+JSV	3.94%	6.64%	1.89%
JFA+JSV	3.49%	5.49%	<b>1.23%</b>
FSH+JFA	3.21%	5.26%	1.34%
FSH+GSV	<b>2.93%</b>	<b>4.06%</b>	1.34%
FSH+JSV	3.21%	4.71%	1.34%
GSV+JSV	3.60%	4.97%	1.34%

## Training Data

### UBM:

- Tel – NIST04, 05, 06 tel speech
- Mic – NIST 05, 06 mic speech

### JFA:

- Matrix  $\mathbf{V}$  – NIST04, 05, 06, Switchboard Phase2, Phase 3, Cellular Parts 2 (300 speaker factors)
- Tel Matrix  $\mathbf{U}$  – NIST04, 05, 06 tel speech (100 Ch factors)
- Mic Matrix  $\mathbf{U}$  – NIST05, 06 mic speech (75 Ch factors)
- Interview Matrix  $\mathbf{U}$  – NIST08 interview speech (75 Ch factors)
- Totally rank of  $\mathbf{U} = 100$  tel +75 mic +75 interview

### Fishvoice:

- Projection matrix  $\mathbf{W}$  – NIST04, 05, 06 tel speech (400 gender-dependent speakers with 8 different utterances)

### NAP:

- Tel – NIST04, 05, 06 tel speech (Corank = 16)
- Mic/interview – NIST05, 06 mic speech, NIST08 interview speech (Corank = 128)

### SVM Imposter-Class:

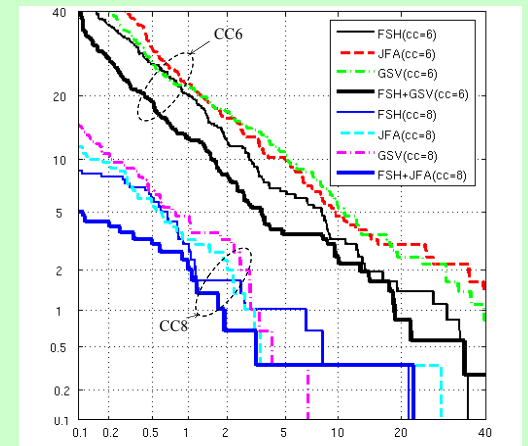
- JSV and JSF – NIST04, 05, 06, 08, Switchboard Cellular Parts 2
- GSV – NIST 06 tel speech and NIST 06 mic speech

### Score normalization:

- Tnorm for JSV, JSF and GSV – NIST04-06 tel/mic speech
- TZnorm for JFA and FSH – NIST04-06 tel/mic speech

EER in % (1<sup>st</sup> line) and Minimum DCF (2<sup>nd</sup> line) for CC1 to CC9.

Com. Cond.	JFA	JSV	JSF	GSV	FSH	Fusion	Rel. Imp.
1	<b>3.88</b> <b>0.63</b>	4.55 0.73	7.51 0.83	4.40 0.67	7.10 0.73	2.69 0.41	31% 35%
2	8.04 0.84	9.11 0.81	13.55 0.89	<b>7.39</b> <b>0.77</b>	11.32 0.86	5.70 0.57	23% 27%
3	<b>4.53</b> <b>0.67</b>	7.59 0.79	11.32 0.89	6.28 0.70	8.26 0.90	2.94 0.51	35% 23%
4	5.79 0.76	7.14 0.68	11.11 0.83	<b>5.58</b> <b>0.68</b>	6.96 0.85	3.59 0.54	36% 20%
5	4.52 <b>0.47</b>	5.73 0.65	5.77 0.60	4.76 0.57	<b>4.09</b> <b>0.54</b>	2.36 0.39	42% 18%
6	7.17 0.82	8.31 0.84	9.14 0.83	7.75 <b>0.79</b>	<b>6.09</b> 0.81	3.87 0.68	36% 14%
7	<b>7.52</b> <b>0.74</b>	8.79 0.91	8.63 0.78	9.15 0.75	8.35 0.85	5.00 0.55	34% 26%
8	2.01 0.46	2.68 0.55	3.69 0.44	2.57 0.48	<b>1.68</b> <b>0.28</b>	1.00 0.22	40% 24%
9	3.45 0.40	<b>2.76</b> 0.46	3.79 0.41	4.13 <b>0.39</b>	4.14 0.49	1.69 0.17	39% 56%



## Conclusions

The HKCUPU system submitted to NIST 2010 SRE is composed of 5 subsystems. The fusion system reduces the EER by 42% and minDCF by 56% when compared with the best individual subsystems. It was also found that the newly proposed FSH subsystem is complementary to JFA and performs significantly better than JFA when its projection matrices were trained by the type of speech that matches the evaluation conditions.