The SYSU System for the Interspeech 2015 Automatic Speaker Verification Spoofing and Countermeasures Challenge*

Shitao Weng[†], Shushan Chen[†], Lei Yu[†], Xuewei Wu[†], Weicheng Cai[‡], Zhi Liu[†], Yiming Zhou[§], Ming Li[‡]

[‡] SYSU-CMU Joint Institute of Engineering, Sun Yat-Sen University, Guangzhou, China

[†] SYSU-CMU Shunde International Joint Research Institute, Guangdong, China

[§] Department of Physics, Zhejiang University, Zhejiang, China

E-mail: wengsht@mail2.sysu.edu.cn

Abstract-Many existing speaker verification systems are reported to be vulnerable against different spoofing attacks, for example speech synthesis, voice conversion, play back, etc. In order to detect these spoofed speech signals as a countermeasure, we propose a score level fusion approach with several different ivector subsystems. We show that the acoustic level Mel-frequency cepstral coefficients (MFCC) features, the phase level modified group delay cepstral coefficients (MGDCC) and the phonetic level phoneme posterior probability (PPP) tandem features are effective for the countermeasure. Furthermore, feature level fusion of these features before i-vector modeling also enhance the performance. A polynomial kernel support vector machine is adopted as the supervised classifier. In order to enhance the generalizability of the countermeasure, we also adopted the cosine similarity and PLDA scoring as one-class classifications methods. By combining the proposed i-vector subsystems with the OpenSMILE baseline which covers the acoustic and prosodic information further improves the final performance. The proposed fusion system achieves 0.29% and 3.26% EER on the development and test set of the database provided by the INTERSPEECH 2015 automatic speaker verification spoofing and countermeasures challenge.

Index Terms: speaker verification, spoofing and countermeasures, i-vector, modified group delay cepstral coefficients, phoneme posterior probability

I. INTRODUCTION

The goal of speaker verification is to automatically verify the claimed speaker identity given a segment of speech. In the past decade, speaker verification has attracted significant research attention with promising results [1]. However, recently it is reported that many existing speaker verification systems are vulnerable against different spoofing attacks, e.g. speech synthesis, voice conversion, play back, etc.[2], [3], [4]

Compared to text independent speaker verification, text dependent speaker verification is more robust against the play back spoofing since the speech content is constrained or pre-defined. Speaker-adapted speech synthesis and voice conversion are the most common spoofing methods that can convert arbitrary text or speech inputs towards the target speaker [2]. To enhance the robustness of speech verification system against spoofing attacks, different countermeasures have been proposed. In [5], higher-level dynamic features and voice quality assessment are used to detect those artificial signals. Furthermore, modified group delay cepstral coefficients (MGDCC) feature has been proposed to distinguish between the original and the spoofed speech signals in the phase domain [6]. This approach is based on the fact that the phase information of synthetic spoofing speech is typically different from the real human articulated speech while the human auditory system is less sensitive to this difference. Apart from acoustic features, prosodic features [7], [8] are also used widely in speech systems. Long term temporal modulation feature derived from magnitude or phase spectrum has also been proposed to detect the synthetic speech [9].

Total variability i-vector modeling has been widely used in speaker verification due to its excellent performance, compact representation and small model size [10], [11]. In this work, we apply the recently proposed generalized i-vector framework [12], [13], [14] with both the acoustic and phonetic features to the countermeasure task.

Figure 1 shows an overview of our anti-spoofing countermeasure system. First, there are several i-vector subsystems using different features, namely the acoustic level Melfrequency cepstral coefficients (MFCC) features, the phase level MGDCC features, the phonetic level phoneme posterior probability (PPP) tandem features [13], [15] and their feature level combinations. Second, we also applied the openS-MILE toolkit [16] to perform the utterance level acoustic and prosodic feature extraction. We believe that the spoofed speech signal may have different prosodic patterns. Third, after the feature normalization, multiple classification methods, e.g. cosine scoring, K-nearest neighbor (KNN), simplified PLDA [17] and Support Vector Machine (SVM), are employed as the back end. Finally, score level fusion is performed to further enhance the overall system performance.

The remainder of the paper is organized as follows. The corpus and the proposed algorithms are explained in Sections II and III, respectively. Experimental results and discussions are presented in Section IV while conclusions are provided in Section V.

^{*}This research was funded in part by the National Natural Science Foundation of China (61401524), Natural Science Foundation of Guangdong Province (2014A030313123), the Fundamental Research Funds for the Central Universities(15lgjc10), CMU-SYSU Collaborative Innovation research Center Foundation(35321.2.8.1011475) and Guangdong Shunde SYSU-CMU Joint Research Institute.



Fig. 1. The system overview

II. CORPUS

The database used to evaluate the proposed methods is based upon a standard dataset of both genuine and spoofed speech. Genuine speech is without significant channel or background noise effect and includes 106 speakers (45 male, 61 female), while spoofed speech is obtained through applying several spoofing algorithms on the genuine speech [18]. The training data set (25 speakers, 3750 genuine utterances and 12635 spoofed utterances) is for model training while the development data set (35 speakers, 3497 genuine utterances and 49875 spoofed utterances) is used to evaluate the system performance and tune the parameters. Finally, the testing data set (46 speakers, 193404 utterances) with unknown types of spoofing attacks is provided to obtain the official submission scores. The details of the database and evaluation protocol are provided in [18].

III. METHODS

From Figure 1, we can see that there are four different features, namely MFCC i-vectors, MFCC-PPP i-vectors, MGDCC-PPP i-vectors and openSMILE feature vectors followed by the same feature normalization, classification and score level fusion pipeline. We first present the proposed features in section III-A. Then section III-B describes the supervised classification and score level fusion methods, respectively.

A. Features

1) The i-vector framework: In the total variability space, there is no distinction between the speaker effects and the channel effects. Rather than separately using the eigenvoice matrix V and the eigenchannel matrix U [19], the total variability space simultaneously captures the speaker and channel variabilities [11]. Given a C component GMM UBM model λ with $\lambda_c = \{p_c, \mu_c, \Sigma_c\}, c = 1, \cdots, C$ and an utterance with a L frame feature sequence $\{y_1, \cdots, y_L\}$, the zero-order and centered first-order Baum-Welch statistics on the UBM are calculated as follows:

$$N_c = \sum_{t=1}^{L} P(c|\mathbf{y}_t, \lambda) \tag{1}$$

$$\mathbf{F_c} = \sum_{t=1}^{L} P(c|\mathbf{y_t}, \lambda)(\mathbf{y_t} - \mu_c)$$
(2)

where $c = 1, \dots, C$ is the GMM component index and $P(c|\mathbf{y}_t, \lambda)$ is the occupancy posterior probability for \mathbf{y}_t on λ_c . The corresponding centered mean supervector $\tilde{\mathbf{F}}$ is generated by concatenating all the $\tilde{\mathbf{F}}_{\mathbf{c}}$ together:

$$\tilde{\mathbf{F}}_{\mathbf{c}} = \frac{\sum_{t=1}^{L} P(c|\mathbf{y}_{\mathbf{t}}, \lambda)(\mathbf{y}_{\mathbf{t}} - \mu_{\mathbf{c}})}{\sum_{t=1}^{L} P(c|\mathbf{y}_{\mathbf{t}}, \lambda)}.$$
(3)

Then the centered mean supervector $\tilde{\mathbf{F}}$ is projected as follows:

$$\tilde{\mathbf{F}} \to \mathbf{T}\mathbf{x},$$
 (4)

where T is a rectangular low rank total variability matrix and x is the so-called i-vector [11].

In the proposed system, 8000 utterances are used to estimate the total variability matrix for a 400 dimensional subspace.

2) The MFCC i-vector: The MFCC i-vector is extracted by the aforementioned i-vector framework with the acoustic level MFCC features. For cepstral feature extraction, a 25ms Hamming window with 10ms shifts was adopted. Each utterance was converted into a sequence of 36-dimensional feature vectors, each consisting of 18 MFCC coefficients and their first order derivatives. We employed the English phoneme recognizer [20] to perform the voice activity detection (VAD) by simply dropping all frames that are decoded as silence or speaker noises.

3) The MFCC-PPP *i*-vector: It is reported in [13], [14] that by combining the phonetic level phoneme posterior probability based tandem features with the acoustic level MFCC features at the feature level, the performances on speaker verification and language identification are significantly enhanced. In this work, the MFCC-PPP i-vector is extracted the same way as in [13] following the generalized i-vector framework. We employed the multilayer perceptron (MLP) based phoneme recognizer [20] with a provided English acoustic model trained on the TIMIT database to perform the phoneme decoding. The GMM model size and the tandem feature dimensionality are 512 and 32, respectively.

4) The MGDCC-PPP i-vector: The MGDCC-PPP i-vector is calculated the same way as the MFCC-PPP i-vector except that here we replace the acoustic level MFCC features with the phase domain MGDCC features. The MGDCC feature is a kind of frame-level feature focusing on the speech phase characteristics. It has been shown that phase domain features are effective for anti-spoofing countermeasures [9]. In order to calculate the MGDCC feature, we need to obtain the modified group delay function phase spectrum (MGDFPS) [21] first.

Given the data x_n of a short time window, the MGDFPS spectrum $\tau_{\rho,\gamma}(\omega)$ is calculated as follows [21]:

$$\tau_{\rho}(\omega) = \frac{X_R(\omega)Y_R(\omega) + Y_I(\omega)X_I(\omega)}{|S(\omega)^{2\rho}|}$$
(5)

$$\tau_{\rho,\gamma}(\omega) = \frac{\tau_{\rho}(\omega)}{|\tau_{\rho}(\omega)|} |\tau_{\rho}(\omega)|^{\gamma}$$
(6)

where $X(\omega)$ and $Y_{(\omega)}$ are the fourier transforms of speech signal x(n) and nx(n); $X_R(\omega)$ and $X_I(\omega)$ are the real and imaginary parts of $X(\omega)$; $Y_R(\omega)$ and $Y_I(\omega)$ are the real and imaginary parts of $Y(\omega)$, respectively. $|S(\omega)|^2$ is calculated by applying a smoothing over $X(\omega)$ [21]. After applying the Mel-frequency filter banks and Discrete Cosine Transform, MGDCC feature is obtained. More details can be found in [9].

5) The OpenSMILE feature vector: The OpenSMILE feature is a 6373 dimensional utterance level feature vector extracted by the OpenSMILE toolkit [16] using the configuration file provided by the 2014 Paralinguistic Challenge [22]. Since various kinds of features, such as MFCC, loudness, auditory spectrum, voicing probability, F0, F0 envelop, jitter, and shimmer, etc., are included, this feature set can capture spoofing information at both the acoustic and prosodic levels. In our system, it served as a baseline as well as a supplement to those i-vector subsystems.

B. Back-end modeling

After feature vectors are extracted, we apply different classification methods for the back-end modeling.

1) The K-nearest neighbor classification (KNN): KNN is a non-parametric multi-class classifier. The utterances in the training set are divided into human set and spoofed set. For each test utterance x_t , K nearest neighboring utterances are found in the training set and the score is calculated based on the class distribution of these K nearest neighbors.

2) The cosine similarity scoring: In our system, a mean vector of all the human utterances in the training data set is calculated. For each test utterance, the score is computed as the cosine similarity between itself and the human class mean vector.

3) PLDA modeling: We first applied the simplified PLDA modeling [17] as the back-end assuming that there are six special speakers (five spoofing channels plus one human channel), each represents a spoofing type or the original genuine speech. Furthermore, we also adopted the two subspace (speaker subspace and spoofing subspace) PLDA presented in [23] to model the i-vectors. The standard log likelihood ratio based hypothesis is emploied for the scoring [17], [23].

4) Support Vector Machine: We formed the anti-spoofing countermeasure as a two class classification task for SVM modeling. The linear kernel LIBLINEAR [24] and its polynomial kernel extension LIBPOLY [25] are adopted as the back-end SVM classifiers and we applied the min/max normalization (range -1 to +1) for each feature dimension on the training, development and test sets with parameters computed only from the training data.

5) Score fusion: We simply employed the weighted summation fusion approach at the score level to further enhance the performance. The fusion weights were tuned on the development data set.

IV. EXPERIMENTAL RESULTS

The results of our four subsystems on the development data are shown in the Table I. We can observe that feature level fusion with PPP feature improves the performance. Compared to the MFCC i-vector subsystem (EER = 6.63%), the EER of MFCC-PPP i-vector subsystem is reduced to 1.06%. On the

other hand, the openSMILE feature outperformed the MFCC i-vector subsystem which might be due to the inclusion of prosodic level information.

Furthermore, to obtain a robust countermeasure system, different backend classification techniques were evaluated. Table II shows the performance on the development data. We used Opensmile with LibSVM as our baseline system and we did not test it on other classifications. Opensmile with SVM classification is the baseline system. Therefore, we didn't try more classifications on this feature. On the other hand, the two stage PLDA classification gives a poor result on MFCC-PPP i-vector feature which is the best feature applied on the other classifications. We didn't consider fusing the result of two stage PLDA classification to the proposed system. Among these six classification methods, LIBPOLY achieves the best performance with an improvement from the performance of baseline system 1.57% to 0.29% EER on the development data. The improvement of LIBPOLY against LIBLINEAR motivated us to further increase the SVM polynomial kernel degree.

As shown in Table III, we simulated unknown spoofing attacks by using four kinds of spoofed utterances in the training and the remaining one in the testing. Although its performance was as good as LIBLINEAR against familiar spoofing attacks (shown in table II), it outperformed LIBLINEAR on the unseen testing data, especially where the unknown attacks were related to speech synthesis (index 3 and 4). This implied that PLDA is more resistive than Liblinear to unseen spoofing attacks. The two stage PLDA only achieved moderate results in Table II which might be because total speakers number in the training data is limited (25) and the speaker subspace may not be orthogonal to the spoofing subspace.

Table IV presents our fusion system results with each individual spoofing condition on the test data. Here S1 to S5 are know attacks and S6 to S10 are unknown attacks. S3 and S4 are synthetic waveform, while S1, S2 and S5 are generated using voice conversion. Our system performed well on all attacks except S10, on which most challenge participants got unsatisfied results.

Finally, our fusion system (system 7) achieved 0.38% and 6.15% EER against known and unknown attacks, respectively.

Methods	EER(LIBPOLY)
MFCC i-vector	6.63
MFCC-PPP i-vector	1.06
MGDCC-PPP i-vector	2.23
OpenSmile	1.57

TABLE I

PERFORMANCE OF SUBSYSTEMS ON THE DEVELOPMENT DATA (LIBPOLY)

V. CONCLUSIONS

This paper presents an anti-spoofing countermeasure system based on a multi-feature and multi-subsystem fusion approach. By fusing the phonetic level phoneme posterior probability tandem features with the acoustic level MFCC features or the phase level MGDCC features, the system performance is significantly enhanced. Combining the proposed i-vector

System	classification method EER	LIBLINEAR	LIBPOLY	COSINE SCORING	KNN	Simplified	two stage
	Feature					PLDA	PLDA
1	MFCC i-vector	8.46	6.63	16.1	9.95	12.01	17.84
2	PPP i-vector	1.72	1.26	3.6	3.4	2.29	-
3	MFCC-PPP i-vector	1.86	1.06	2.86	2.46	1.89	10.18
4	MGDCC-MFCC-PPP i-vector	2.97	2.06	6.52	3.43	3.95	17.79
5	OPENSmile	2.03	1.57	-	-	-	-
6	Fusion 1+2+3+4	-	-	1.63	1.37	1.09	-
7	Fusion 1+2+3+4+5	0.54	0.29	-	-	-	-

TABLE II

PERFORMANCE OF THE PROPOSED METHODS ON THE DEVELOPMENT DATA

	S1	S2	\$3	S4	S5	S6	S 7	S8	S9	S10	Average
EER (Fusion 1+2+3+4+5-LIBPOLY)	0.1137	1.0332	0.0482	0.0412	0.6614	0.7112	0.2297	0.0108	0.1336	29.6649	3.265

TABLE IV

PERFORMANCE OF THE FUSION SYSTEMS WITH DIFFERENT SPOOFING CONDITIONS ON THE TESTING DATA

train set	test set	PLDA	LIBLINEAR
human+spoof[2,3,4,5]	human+spoof[1]	3.57	3.4
human+spoof[1,3,4,5]	human+spoof[2]	4.8	7.69
human+spoof[1,2,4,5]	human+spoof[3]	0.2	0.71
human+spoof[1,2,3,5]	human+spoof[4]	0.2	0.66
human+spoof[1,2,3,4]	human+spoof[5]	4.49	11.81

TABLE III

PERFORMANCE (EER) OF THE LIBLINEAR AND THE SIMPLIFIED PLDA BACKENDS ON THE UNKNOWN SPOOFING TESTING CONDITIONS

subsystems with the OpenSMILE baseline which covers the acoustic and prosodic level information further improves the final performance. For the back-end modeling, two classes support vector machine outperforms the one class cosine similarity or PLDA scoring on the development data where the spoofing attack types are known. The one class scoring method achieves more robust performance on the unseen testing data where the spoofing conditions are unknown.

REFERENCES

- T. Kinnunena 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] Z. Wu, N. W. D. Evans, T. Kinnunen, J. Yamagishi, F. Alegre, and H. Li, "Spoofing and countermeasures for speaker verification: A survey," *Speech Communication*, vol. 66, pp. 130–153, 2015.
- [3] Z. Wu, T. Virtanen, T. Kinnunen, E. Chng, and H. Li, "Exemplar-based unit selection for voice conversion utilizing temporal information," in *proceedings of INTERSPEECH*, 2013, pp. 3057–3061.
 [4] T. Toda, A. W. Black, and K. Tokuda, "Voice conversion based on
- [4] T. Toda, A. W. Black, and K. Tokuda, "Voice conversion based on maximum-likelihood estimation of spectral parameter trajectory," *IEEE Transactions on Audio, Speech & Language Processing*, vol. 15, no. 8, pp. 2222–2235, 2007.
- [5] F. Alegre, R. Vipperla, and N. Evans, "Spoofing countermeasures for the protection of automatic speaker recognition systems against attacks with artificial signals," in *proceedings of INTERSPEECH*, 2012.
- [6] Z. Wu, C. E. Siong, and H. Li, "Detecting converted speech and natural speech for anti-spoofing attack in speaker recognition." in *proceedings* of *INTERSPEECH*, 2012.
- [7] E. Shriberg, L. Ferrer, S. S. Kajarekar, A. Venkataraman, and A. Stolcke, "Modeling prosodic feature sequences for speaker recognition," *Speech Communication*, vol. 46, no. 3-4, pp. 455–472, 2005. [Online]. Available: http://dx.doi.org/10.1016/j.specom.2005.02.018
- [8] C. C. Leung, M. Ferras, C. Barras, and J. L. Gauvain, "Comparing prosodic models for speaker recognition," in *INTERSPEECH 2008*, 9th Annual Conference of the International Speech Communication Association, Brisbane, Australia, September 22-26, 2008, 2008, pp. 1945–1948.
- [9] Z. Wu, X. Xiao, E. Chng, and H. Li, "Synthetic speech detection using temporal modulation feature," in *Proceedings of ICASSP*, 2013, pp. 7234–7238.

- [10] N. Dehak, P. Torres-Carrasquillo, D. Reynolds, and R. Dehak, "Language recognition via i-vectors and dimensionality reduction," in *Proc. INTERSPEECH*, 2011, pp. 857–860.
- [11] N. Dehak, P. Kenny, R. Dehak, P. Dumouchel, and P. Ouellet, "Front-end factor analysis for speaker verification," *IEEE Transactions on Audio*, *Speech, and Language Processing*, vol. 19, no. 4, pp. 788–798, 2011.
- [12] Y. Lei, N. Scheffer, L. Ferrer, and M. McLaren, "A novel scheme for speaker recognition using a phonetically-aware deep neural network," in *Proc. ICASSP*, 2014.
- [13] M. Li and W. Liu, "Speaker verification and spoken language identification using a generalized i-vector framework with phonetic tokenizations and tandem features," in *Proc. INTERSPEECH*, 2014.
- [14] L. D'Haro, R. Cordoba, C. Salamea, and J. Echeverry, "Extended phone log-likelihood ratio features and acoustic-based i-vectors for language recognition," in *Proc. ICASSP.* IEEE, 2014, pp. 5379–5383.
- [15] M. Li, "Automatic recognition of speaker physical load using posterior probability based features from acoustic and phonetic tokens," in *Proc. INTERSPEECH*, 2014.
- [16] F. Eyben, M. Wöllmer, and B. Schuller, "Opensmile: the munich versatile and fast open-source audio feature extractor," in *Proceedings* of the international conference on Multimedia. ACM, 2010, pp. 1459– 1462.
- [17] D. Garcia-Romero and C. Y. Espy-Wilson, "Analysis of i-vector length normalization in speaker recognition systems," in *Proceedings of IN-TERSPEECH*, 2011, pp. 249–252.
- [18] Z. Wu, T. Kinnunen, N. Evans, and J. Yamagishi, "Asvspoof 2015: Automatic speaker verification spoofing and countermeasures challenge evaluation plan." [Online]. Available: http://www.spoofingchallenge.org/
- [19] P. Kenny, G. Boulianne, P. Ouellet, and P. Dumouchel, "Joint factor analysis versus eigenchannels in speaker recognition," *IEEE Transactions on Audio, Speech, and Language Processing*, vol. 15, no. 4, pp. 1435–1447, 2007.
- [20] P. Schwarz, P. Matejka, and J. Cernocky, "Hierarchical structures of neural networks for phoneme," in *Proc. ICASSP*, 2006, pp. 325–328, software available at http://speech.fit.vutbr.cz/software/phoneme-recognizer-based-long-temporal-context.
 [21] D. Zhu and K. K. Paliwal, "Product of power spectrum and group delay
- [21] D. Zhu and K. K. Paliwal, "Product of power spectrum and group delay function for speech recognition," in *Proceedings of ICASSP*, 2004, pp. 125–128.
- [22] B. Schuller, S. Steidl, A. Batliner, F. Epps, J.and Eyben, F. Ringeval, E. Marchi, and Y. Zhang, "The interspeech 2014 computational paralinguistics challenge: Cognitive & physical load," in *Proc. INTERSPEECH*, 2014.
- [23] S. Prince and J. Elder, "Probabilistic linear discriminant analysis for inferences about identity," in *Proc. ICCV*, 2007, pp. 1–8.
- [24] R.-E. Fan, K.-W. Chang, C.-J. Hsieh, X.-R. Wang, and C.-J. Lin, "Liblinear: A library for large linear classification," *Journal of machine learning research*, vol. 9, pp. 1871–1874, 2008.
- [25] Y.-W. Chang, C.-J. Hsieh, K.-W. Chang, M. Ringgaard, and C.-J. Lin, "Training and Testing Low-degree Polynomial Data Mappings via Linear SVM," *Journal of Machine Learning Research*, vol. 11, pp. 1471–1490, 2010.