



ELSEVIER

Available online at [www.sciencedirect.com](http://www.sciencedirect.com)

ScienceDirect

Computer Speech &amp; Language xxx (2017) xxx-xxx

[www.elsevier.com/locate/csl](http://www.elsevier.com/locate/csl)

# Synthetic speech detection using fundamental frequency variation and spectral features<sup>☆</sup>

Monisankha Pal<sup>\*,a</sup>, Dipjyoti Paul<sup>a</sup>, Goutam Saha<sup>a</sup>

Department of Electronics and Electrical Communication Engineering, Indian Institute of Technology Kharagpur, Kharagpur 721 302, India

Received 18 June 2017; received in revised form 31 August 2017; accepted 6 October 2017

Available online xxx

## Abstract

Recent works on the vulnerability of automatic speaker verification (ASV) systems confirm that malicious spoofing attacks using synthetic speech can provoke significant increase in false acceptance rate. A reliable detection of synthetic speech is key to develop countermeasure for synthetic speech based spoofing attacks. In this paper, we targeted that by focusing on three major types of artifacts related to magnitude, phase and pitch variation, which are introduced during the generation of synthetic speech. We proposed a new approach to detect synthetic speech using score-level fusion of front-end features namely, constant Q cepstral coefficients (CQCCs), all-pole group delay function (APGDF) and fundamental frequency variation (FFV). CQCC and APGDF were individually used earlier for spoofing detection task and yielded the best performance among magnitude and phase spectrum related features, respectively. The novel FFV feature introduced in this paper to extract pitch variation at frame-level, provides complementary information to CQCC and APGDF. Experimental results show that the proposed approach produces the best stand-alone spoofing detection performance using Gaussian mixture model (GMM) based classifier on ASVspoof 2015 evaluation dataset. An overall equal error rate of 0.05% with a relative performance improvement of 76.19% over the next best-reported results is obtained using the proposed method. In addition to outperforming all existing baseline features for both known and unknown attacks, the proposed feature combination yields superior performance for ASV system (GMM with universal background model/*i*-vector) integrated with countermeasure framework. Further, the proposed method is found to have relatively better generalization ability when either one or both of copy-synthesized data and limited spoofing data are available a priori in the training pool.

© 2017 Elsevier Ltd. All rights reserved.

**Keywords:** All-pole group delay function (APGDF); Anti-spoofing; Constant Q cepstral coefficient (CQCC); Fundamental frequency variation (FFV); Score-level fusion; Spoofing attack

## 1. Introduction

Automatic speaker verification (ASV) systems that accept or reject an identity claim have a wide range of applications in banking, forensics, voice mail, etc. (Kinnunen and Li, 2010; Pal and Saha, 2015). However, the main concern with the deployment of ASV systems is their vulnerability towards *spoofing attacks* in which a fraudster tries to

<sup>☆</sup> This paper has been recommended for acceptance by Prof. R. K. Moore.

\* Corresponding Author.

E-mail address: [monisankhapal@iitkgp.ac.in](mailto:monisankhapal@iitkgp.ac.in) (M. Pal), [dipjyotipaul@ece.iitkgp.ernet.in](mailto:dipjyotipaul@ece.iitkgp.ernet.in) (D. Paul), [gsaha@ece.iitkgp.ernet.in](mailto:gsaha@ece.iitkgp.ernet.in) (G. Saha).

5 masquerade an enrolled person's voice to get illegitimate acceptance. The major forms of spoofing attacks are *imper-*  
6 *sonation* (Hautamäki et al., 2015), *replay* (Villalba and Lleida, 2011), *speaker adapted speech synthesis* (Moulines  
7 and Charpentier, 1990) and *voice conversion* (Stylianou et al., 1998). Recent advancements in computer-assisted  
8 speech synthesis (SS) and voice conversion (VC) technology, with availability of related open source software,  
9 make SS and VC the most potent means of spoofing attacks to swindle ASV systems. To mitigate this, appropriate  
10 *countermeasure* (CM) to discriminate natural speech from the synthetic speech is essential.

11 A popular approach for spoofing detection is the use of front-end features that targets capturing the artifacts intro-  
12 duced during speech manipulation. Modulation features extracted from magnitude or phase spectrum, which carry  
13 long-term temporal information were incorporated to distinguish HMM-based synthetic speech from natural speech  
14 (Wu et al., 2013). Based on human speech perception and the fact that phase information is usually lost during syn-  
15 thesis in VC, phase-based features like cosine-normalized phase (CosPhase) and modified group delay function  
16 (MGDF) were proposed to detect VC spoofing as stand-alone system (Wu et al., 2012b) and with ASV framework  
17 (Wu et al., 2012a). The use of relative phase shift (RPS) features for reliable detection of hidden Markov model  
18 (HMM) based text-to-speech (TTS) synthesized speech was introduced in De Leon et al. (2012a). Another phase-  
19 based feature called mel regularized RPS was also explored to detect TTS synthesized speech (Sanchez et al.,  
20 2015b). A synthetic speech detector using prosodic features like pitch pattern statistics from image analysis of pitch  
21 patterns was proposed in De Leon et al. (2012b). A new countermeasure based on spectro-temporal information  
22 from local binary patterns was derived in Alegre et al. (2013), which is less reliant on prior knowledge and provides  
23 robust protection from VC as well as SS attacks. Countermeasure using high-level dynamic features and voice qual-  
24 ity assessment was employed in Alegre et al. (2012). Most of these hand-crafted discriminative features used prior  
25 spoofing data to distinguish natural and synthetic utterances drawn from a closed set database. However, develop-  
26 ment of generalized countermeasure to act on unseen spoofing data type is gaining practical importance. Apart from  
27 focusing on front-end features, *anti-spoofing* using i-vector based representation of speech utterances as back-end  
28 was investigated in Sizov et al. (2015). This back-end generative model is promising to provide generalized counter-  
29 measure.

30 To address the issues arising from using closed set databases, the research community developed ASVspooF 2015  
31 (Wu et al., 2015), which contains spoofed data from a diverse range of spoofing attacks. The winning system (Patel  
32 and Patil, 2017a) of this challenge used a combination of standard mel frequency cepstral coefficients (MFCCs) and  
33 cochlear filter cepstral coefficients (CFCCs) with the change in instantaneous frequency (IF). A significant number  
34 of countermeasures have been proposed till date on this challenge data. Among them, MGD-phase features (Alam  
35 et al., 2015; Liu et al., 2015), RPS features (Sanchez et al., 2015a; Wang et al., 2015), discriminatory sub-band fea-  
36 tures (Sriskandaraja et al., 2016), wavelet-based features (Novoselov et al., 2016), linear prediction features (Alam  
37 et al., 2015; Janicki, 2015) were explored for spoofing detection task. A comparative study showing the efficacy of  
38 various short-term power spectrum and phase features, dynamic and complementary features for anti-spoofing was  
39 presented in Sahidullah et al. (2015). In this study, linear frequency cepstral coefficient (LFCC) and inverted mel fre-  
40 quency cepstral coefficient features emerged as excellent discriminative features for spoofing detection. Short-term  
41 spectral features like *constant Q cepstral coefficients* (CQCCs) exhibited best results on ASVspooF 2015 data as  
42 found in Todisco et al. (2017). Recently, scattering cepstral coefficients (SCCs) which are similar to CQCCs were  
43 proposed for stand-alone spoofing detection (Sriskandaraja et al., 2017). The effectiveness of cepstral features for  
44 spoofing detection as stand-alone and integrated with ASV using inverted frequency warping scale and overlapped  
45 block transformation was demonstrated in Paul et al. (2017a). At classifier level, the usage of deep neural network  
46 (DNN) (Chen et al., 2015; Soni et al., 2016), DNN and support vector machine (Villalba et al., 2015), i-vectors  
47 (Weng et al., 2015) were used on open set ASVspooF data.

48 The standard approach for spoofing detection captures artifact traces introduced while doing speech manipulation  
49 by SS or VC. This is done by using efficient feature extraction techniques as front-end and standard classifiers at the  
50 back-end. The work in Hanilçi et al. (2015) suggests that more reliable spoofing detection can be achieved by dis-  
51 criminative features than complex classifiers. However, most of the countermeasures in literature focused on a par-  
52 ticular type of artifacts (phase, pitch contour, high frequency spectral, temporal). ASVspooF 2015 challenge results  
53 also showed that most of the submitted systems provided poorer performance for unknown spoofed data types (Wu  
54 et al., 2015). It suggests that instead of a single feature, a bank of features may prove to be useful for spoofing detec-  
55 tion. Therefore, the aim of this work is to propose an efficient feature combination that can capture possible sources  
56 of artifacts and can easily distinguish natural speech from synthetic speech for both known and unknown spoofing

Download English Version:

<https://daneshyari.com/en/article/6951503>

Download Persian Version:

<https://daneshyari.com/article/6951503>

[Daneshyari.com](https://daneshyari.com)