
Speaker identification for aeronautical communications

Sara SEKKATE
LIM@II FSTM
sarasekkate@gmail.com

Mohammed KHALIL
LIM@II FSTM
medkhalil87@gmail.com

Abdellah ADIB
LIM@II FSTM
adib@fstm.ac.ma

Abstract

1 In this paper, we aim to fulfill this goal by developing a Speaker Identification
2 System (SIS) for future aeronautical communications systems. Furthermore, we
3 present a novel feature extraction scheme based on multi-resolution analysis.
4 The proposed system, called i-SMFCC uses Mel Frequency Cepstral Coefficients
5 (MFCC) features of Stationary Wavelet Transform (SWT) sub-bands. The extracted
6 features are modeled using the i-vector approach and Support Vector Machines
7 (SVM) are adopted as a back-end classifier. The performance of the proposed SIS
8 is evaluated using two publicly available databases. Comparison of the proposed
9 approach with the baseline MFCC feature extraction shows the feasibility and the
10 robustness of the proposed method.

11 1 Introduction

12 Speaker recognition is a well-established research problem and has found use in many applications,
13 including voice authenticated bank transaction, access control, and prison call monitoring [1]. How-
14 ever, research has been rarely devoted to the integration of Automatic Speaker Recognition (ASR) in
15 the aeronautical industry [2].

16 In Air Traffic Control (ATC), voice communication serves as the main media for delivering instructions
17 and important information between pilots and controllers [3]. The road to improving safety in ATC,
18 therefore, definitely passes through improving the air-ground communications safety. By far, the
19 most prominent issue surrounding these communications is the heightened risk of callsign confusion
20 [4]. As a consequence, it is possible for the pilot to accept clearances meant for others, leading to
21 wrong subsequent actions and incidents with a high potential to cause death [5].

22 In this paper, we aim to develop a closed-set¹ text-independent² speaker identification³ for the
23 En-route⁴ airspace that could be used to prevent call-sign confusion and hence increase flight safety.

24 2 The proposed approach

25 The proposed system is summarized as following:

¹Closed-set systems suppose that the unknown speaker to be identified is known a priori to be one of the registered speakers set.

²Text-independent mode imposed no restrictions on spoken phrases, and hence the speaker can utter any word in order to be recognized.

³The goal is to identify an input speech by selecting one model from the previously enrolled speaker models.

⁴En-route scenario describes the state where the aircraft is airborne and communicates with the control tower (air-ground communication) or with an other airplane (air-air communication).

- 26 1. First, for each speech sample, the signal is decomposed into 9 sub-bands using SWT⁵ to
27 achieve a series of approximation and detail coefficients.
- 28 2. Then, MFCC are computed from each sub-band [6]. Here, the first 13 coefficients derived
29 from a 20-channel mel-scaled filterbank are extracted from speech frames of 25ms with
30 a frame shift of 10 ms, removing the first one because it carries less speaker specific
31 information. In addition to static MFCC features, the log-energy as well as first and second
32 derivatives were also included to produce a feature vector of 39 elements.
- 33 3. Concatenate all the sub-band features to produce a final feature vector, denoted as SMFCC.
- 34 4. Repeat step 1 to step 3 for each speech sample to create a feature matrix that will be fed into
35 the i-vector modeling framework [7].
- 36 5. The resulting feature vectors are then assembled, modeled using the i-vector approach and
37 fed to the SVM classifier.

38 3 Experimental results

39 The performances of the proposed system are evaluated on two publicly available datasets. The first
40 one consists of 10 speakers (4 females and 6 males) from ATCOSIM speech corpus [8]. The second
41 one is a set of 455 speakers (41 females and 414 males) from Voxforge database [9].

42 From the results listed in Table 1 for ATCOSIM database, we observe an increase of 5 to 11 percent in
43 the accuracy at 5dB with the highest accuracy of 49% reached for SMFCC features without i-vector
44 modeling. We can also see that in general, the performance of i-MFCC is lower than MFCC at all
45 SNR ranges, while i-SMFCC outperforms SMFCC beyond 10dB. On the other hand, Table 2 reports
46 the results obtained using Voxforge database, where the advantage of using SMFCC is again obvious
47 compared to baseline MFCC. Except for 5dB, although there is an improvement compared to MFCC,
48 the achieved performance remains poor.

Table 1: Accuracy (%) of the proposed method in noisy environments using ATCOSIM speech corpus

Transmission Channel	SNR (dB)	ATCOSIM			
		MFCC	SMFCC	i-MFCC	i-SMFCC
AWGN	5	45.67	49	14.67	25
	10	84.33	92.67	81.67	96.33
	15	84.33	92.67	81.67	96.33
En-route	Infinite	84.33	92.67	81.67	96.33
En-route + AWGN	5	37	36	10.33	15.33
	10	84	89	80.33	93
	15	84.67	93	81.67	96.33

Table 2: Accuracy (%) of the proposed method in noisy environments using Voxforge speech corpus

Transmission Channel	SNR (dB)	Voxforge			
		MFCC	SMFCC	i-MFCC	i-SMFCC
AWGN	5	2.12	10.84	3	5.71
	10	24.47	71.79	79.63	92.45
	15	24.47	71.87	79.63	92.45
En-route	Infinite	24.47	71.87	79.56	92.45
En-route + AWGN	5	1.39	8.64	1.47	2.56
	10	22.12	66.81	75.53	84.40
	15	24.40	71.72	79.56	92.45

49 References

- 50 [1] Zhao, X. & Wang, Y. & Wang, D. (2014): Robust speaker identification in noisy and reverberant
51 conditions. *IEEE/ACM Transactions on Audio, Speech, and Language Processing* 22(4), 836-845.

⁵The main reason of using SWT is the fact that it is time-invariant transform as compared to DWT.

- 52 [2] Hofbauer, K. & Hering, H. & Kubin, G. (2005): Speech watermarking for the VHF radio channel.
53 *4th EUROCONTROL Innovative Research Workshop*
- 54 [3] Haas, E. (2002) Aeronautical channel modeling. *IEEE Trans. on Vehicular Technology* 51(2),
55 254-264
- 56 [4] Delpech, E. & Laignelet, M. & Pimm, C. & Raynal, C. & Trzos, M. & Arnold, A. & Pronto,
57 D. (2018) A Real-life, French-accented Corpus of Air Traffic Control Communications. *Language*
58 *Resources and Evaluation Conference (LREC), Miyazaki, Japan*
- 59 [5] Chen, S. & Kopald, H. & Chong, R. & Wei, Y. & Levonian, Z. (2017) Read Back Error Detection
60 using Automatic Speech Recognition *Twelfth USA/Europe Air Traffic Management Research and*
61 *Development Seminar*
- 62 [6] Nagaraja, B.G. & Jayanna, H.S. (2012) Multilingual Speaker Identification with the Constraint
63 of Limited Data Using Multitaper MFCC. *Recent Trends in Computer Networks and Distributed*
64 *Systems Security pp. 127-134*
- 65 [7] Kheder, W.B. & Matrouf, D. & Bousquet, P.M. & Bonastre, J.F. & Ajili, M. (2017) Fast i-vector
66 denoising using map estimation and a noise distributions database for robust speaker recognition.
67 *Computer Speech and Language pp.104-122.*
- 68 [8] Hofbauer, K. & Petrik, S. & Hering, H. (2008) The ATCOSIM corpus of non-prompted clean air
69 traffic control speech. *Proceedings of the Sixth International Conference on Language Resources and*
70 *Evaluation, Marrakech, Morocco*
- 71 [9] Voxforge database. Tech. rep. URL <http://voxforge.org>