

## Joint Single-Channel Speech Separation and Speaker Identification

Mowlaee, Pejman; Saeidi, Rahim ; Tan, Zheng-Hua; Christensen, Mads Græsbøll; Fränti, Pasi ; Jensen, Søren Holdt

*Published in:*

I E E E International Conference on Acoustics, Speech and Signal Processing. Proceedings

*DOI (link to publication from Publisher):*

[10.1109/ICASSP.2010.5495619](https://doi.org/10.1109/ICASSP.2010.5495619)

*Publication date:*

2010

*Document Version*

Early version, also known as pre-print

[Link to publication from Aalborg University](#)

*Citation for published version (APA):*

Mowlaee, P., Saeidi, R., Tan, Z.-H., Christensen, M. G., Fränti, P., & Jensen, S. H. (2010). Joint Single-Channel Speech Separation and Speaker Identification. *I E E E International Conference on Acoustics, Speech and Signal Processing. Proceedings, 2010*, 4430 - 4433 . <https://doi.org/10.1109/ICASSP.2010.5495619>

### General rights

Copyright and moral rights for the publications made accessible in the public portal are retained by the authors and/or other copyright owners and it is a condition of accessing publications that users recognise and abide by the legal requirements associated with these rights.

- Users may download and print one copy of any publication from the public portal for the purpose of private study or research.
- You may not further distribute the material or use it for any profit-making activity or commercial gain
- You may freely distribute the URL identifying the publication in the public portal -

### Take down policy

If you believe that this document breaches copyright please contact us at [vbn@aub.aau.dk](mailto:vbn@aub.aau.dk) providing details, and we will remove access to the work immediately and investigate your claim.

# Signal-to-Signal Ratio Independent Speaker Identification for Co-Channel Speech Signals

Rahim Saeidi\*, Pejman Mowlae<sup>†</sup>, Tomi Kinnunen\*, Zheng-Hua Tan<sup>†</sup>, Mads Græsbøll Christensen<sup>‡</sup>, Søren Holdt Jensen<sup>†</sup>, and Pasi Franti\*

\*School of Computing, University of Eastern Finland, Joensuu, Finland  
Email: {rahim.saeidi, tomi.kinnunen, pasi.franti}@uef.fi

<sup>†</sup>Department of Electronic Systems, Aalborg University, Denmark  
Email: {pmb,zt,shj}@es.aau.dk

<sup>‡</sup>Department of Architecture Design and Media Technology  
Aalborg University, Denmark  
Email: mgc@imi.aau.dk

**Abstract**—In this paper, we consider speaker identification for the co-channel scenario in which speech mixture from speakers is recorded by one microphone only. The goal is to identify both of the speakers from their mixed signal. High recognition accuracies have already been reported when an accurately estimated signal-to-signal ratio (SSR) is available. In this paper, we approach the problem without estimating SSR. We show that a simple method based on fusion of adapted Gaussian mixture models and Kullback-Leibler divergence calculated between models, achieves an accuracy of 97% and 93% when the two target speakers enlisted as three and two most probable speakers, respectively.

**Keywords**—Speaker Identification; GMM; MAP adaptation; co-channel speech;

## I. INTRODUCTION

Speaker identification (SID) is the task of recognizing one's identity based on observed speech signal [1]. Typical speaker identification systems consist of short-term spectral feature extractor (front-end) and a pattern matching module (back-end). In traditional SID, the basic assumption is that only one target speaker exists in the given signal whereas in *co-channel* SID, the task is to identify two target speakers in one given mixture. Distinct from the so-called *summed channel* speaker recognition task [2], where only one speaker is talking most of the time, in the co-channel SID problem, both speakers talk simultaneously. Research on co-channel speaker identification has been done for more than one decade [3], yet the problem remains largely unsolved.

Most of the current *single-channel speech separation* (SCSS) systems use a model-based SID module, known as *Iroquois* [4] to identify the speakers in a mixed signal. The goal of an SCSS system is to estimate the unknown speaker signals according to their observed mixture. Interaction of

the SID and speech separation modules can be managed in a closed loop to increase the overall performance [5]. Recognition accuracy as high as 98% has been reported for *Iroquois* in [6] which makes it as a first choice to be included in SCSS systems [7]. The database in [6] is provided for speech separation challenge and consists of 2 seconds of small vocabulary speech for 34 speakers. In the *Iroquois* system, a short list of the most likely speakers are produced based on the frames of the mixed signal that are dominated by one speaker. This short-list is then passed to a *max-based EM algorithm* to find the signal-to-signal ratio (SSR) and two speakers identity with an exhaustive search on codebooks created for speech synthesis [4].

The SSR estimation in *Iroquois* system is based on finding the most likely combination of speakers codebooks to produce the current speech frame, where in text-independent case gets more challenging compared to the database in [6]. Although the SSR can be continuous and time-varying over a recording in realistic conditions, in database presented in [6] and in this study the discrete SSR levels of  $\{-9, -6, -3, 0, 3, 6\}$  dB are considered. Furthermore, in real-time applications of SCSS and in forensic applications it is necessary to have a *fast* and *accurate* system to identify the underlying sources in mixed signal without SSR estimation required.

To this end, in this paper, we propose an SSR-independent SID module for co-channel speech. More specifically, we examine different frame-level likelihood scores and model level distances to solve the problem and propose a combination of the most successful ones to compare the accuracy with respect to *Iroquois*. Since the proposed system is SSR-independent and tuned on 8 kHz speech, it is believed that it could be an alternative approach for the SID in SCSS and useful for telephony data found, for instance, in forensic applications.

The work of R. Saeidi was supported by a scholarship from the Finnish Foundation for Technology Promotion (TES). The work of P. Mowlae is supported by the Marie Curie EST-SIGNAL Fellowship, contract no. MEST-CT-2005-021175. The work of T. Kinnunen was supported by the Academy of Finland (project no 132129).

## II. SPEAKER RECOGNITION APPROACH

We use two main approaches for speaker recognition: frame-level log-likelihood calculation for a given mixed signal against a speaker GMM and between-models distance of a GMM model trained on mixed signal to speaker GMMs.

### A. Frame-level likelihood scores

From the frame-level likelihood estimation originally defined for the *Iroquois* system in [4], [8] and which aims at determining the frames where only one speaker exists, we derive three different scores defined at the end of this section. A maximum likelihood (ML) trained GMM has been used in [4]; however, maximum *a posteriori* (MAP) derived GMMs [9] are more accurate in speaker verification and we follow this latter approach. Let  $\lambda$  denote speaker GMM. The likelihood function is defined as,

$$\ell(\mathbf{x}) = p(\mathbf{x}|\lambda) = \sum_{m=1}^M w_m p_m(\mathbf{x}). \quad (1)$$

The density is a weighted linear combination of  $M$  unimodal Gaussian densities  $p_m(\mathbf{x})$ , where  $p_m(\mathbf{x}) \sim \mathcal{N}(\mathbf{x}; \mu_m, \Sigma_m)$  and the mixture weights  $w_m$  further satisfy the constraints  $\sum_{m=1}^M w_m = 1$  and  $w_m \geq 0$ . Speaker-dependent GMMs are adapted from universal background model (UBM) [9]. The UBM is a GMM trained on a pool of feature vectors extracted from as many speakers as possible to serve as *a priori* information for feature distribution. GMM means are the only parameters updated and weights and covariances are copied directly from UBM to GMMs.

### B. Model distance scores

We define  $\lambda_{ig}$  as the SSR-dependent model for  $i$ th speaker at SSR level  $g$ . Another approach to measure similarity of a speech segment with a speaker model ( $\lambda_i$ ) is to make a model from the test utterance with MAP adaptation ( $\lambda_e$ ) and calculate the distance between  $\lambda_e$  and the speaker model. We use the *Kullback-Leibler divergence* (KLD) as a distance measure between the two probability distributions. Since this distance cannot be directly evaluated for GMMs, we use the upper bound of KLD which has successfully been applied to speaker verification [10]:

$$\text{KLD}_i = \frac{1}{2} \sum_{g=1}^G \sum_{m=1}^M w_m (\mu_{me} - \mu_{mig})^T \Sigma_m^{-1} (\mu_{me} - \mu_{mig}). \quad (2)$$

Here  $G$  ranges in a set of SSR levels,  $\mu_{me}$  is the  $m$ th mean vector in  $\lambda_e$  and  $\mu_{mig}$  is the  $m$ th mean vector in  $\lambda_{ig}$ , whereas  $w_m$  and  $\Sigma_m$  are the weights and the covariances of the UBM, respectively. An alternative approach to measure the distortion between GMMs is *approximate cross entropy* (ACE) [11]. As shown in [11], assuming infinite number of test utterance feature vectors, log-likelihood for a given  $\lambda_i$

equals to negative cross entropy between  $\lambda_e$  and  $\lambda_i$ . It can be approximated as follows:

$$\begin{aligned} \text{ACE}_i = & \sum_{g=1}^G \sum_{m=1}^M w_m \max_n \left[ \log w_n \right. \\ & - \frac{1}{2} (\mu_{me} - \mu_{nig})^T \Sigma_n^{-1} (\mu_{me} - \mu_{nig}) \\ & \left. - \frac{1}{2} \log |\Sigma_n| - \frac{D}{2} (1 + \log 2\pi + \frac{1}{Tw_m + r}) \right], \end{aligned} \quad (3)$$

where  $T$  is the total number of frames for training  $\lambda_e$ ,  $D$  is features dimension and  $r$  is a relevance factor that controls compromise between UBM statistics and adaptation data in GMM adaptation [9]. The value  $r = 0$  corresponds to barely standing on adaptation data.

### C. Proposed method

In this work, we train the UBM ( $\lambda_{UBM}$ ) using digitally mixed speech signals at different SSR levels formed by different speakers. Moreover, we train each target speaker  $i$ , the set of gain-dependent models  $\lambda_{ig}$  that are adapted from the UBM based on  $i$ th speaker speech files corrupted by other speakers signal at SSR level  $g$ . Using SSR-based speaker models, the system captures speaker-dependent information when it is contaminated by other speakers data. This is similar to the idea of having an SSR-based bias in GMM parameters in [4], however, it has the major difference that we build separate GMMs for each SSR level based on the UBM. It enables the system to function independent of the SSR level.

For a feature vector extracted from a speech segment at time instance  $t$ , and denoted by  $\mathbf{x}_t$ , frame level score for speaker  $i$  is defined as,

$$s_{it} = \frac{1}{G} \sum_{g=1}^G \log[p(\mathbf{x}_t|\lambda_{ig})] - \log[p(\mathbf{x}_t|\lambda_{UBM})], \quad (4)$$

We average over all SSR levels to be independent of the underlying SSR in the given signal and normalize all speakers scores at time instance  $t$  with the corresponding UBM score. To emphasize dominant speaker score in a frame, the score in (4) is further normalized by  $s'_{it} = s_{it}/\sigma_t$ , where  $\sigma_t$  is standard deviation of all speakers scores for the frame  $t$ . To sum up, we consider five different scores for a speaker:

**NWF:** *number of winning frames*, where speaker  $i$  is the most probable speaker in that frame,  $NWF_i = \sum_t \varphi(s'_{it})$  where  $\varphi(s'_{it}) = 1$  for  $i = \arg \max_j s'_{jt}$  and 0 otherwise.

**NCF:** *number of confident frames* for speaker  $i$  where  $s'_{it}$  is above threshold  $\alpha$ :  $NCF_i = \sum_t \psi(s'_{it})$  where  $\psi(s'_{it}) = 1$  for  $s'_{it} > \alpha$  and 0 otherwise.

**LL:** *Log-likelihood mean* for which  $s'_{it}$  is above threshold  $\alpha$ :  $LL_i = (1/NCF_i) \sum_t \psi(s'_{it}) s'_{it}$ .

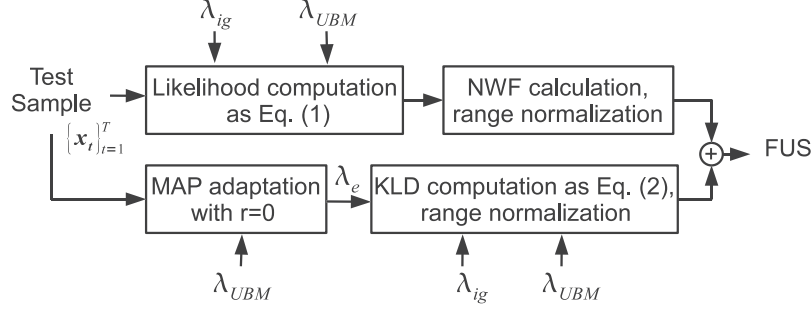


Figure 1. Proposed SID module is a combination of frame level likelihood score and model level distance:  $FUS = 0.54NWF + 0.46KLD$ .

**KLD:** *Kullback-Leibler divergence* between  $\lambda_e$  and a set of models  $\lambda_{ig}$ , computed using (2).

**ACE:** *approximate cross entropy* between  $\lambda_e$  and a set of models  $\lambda_{ig}$ , computed using (3).

As it is common in speaker recognition, to enable using benefits from different recognizers, we considered the fusion of the scores. We used an approximate brute-force search to find the optimal weights for score fusion. It should be mentioned that we normalized (and reverted for KLD) the range of scores from different recognizers before fusion. A block diagram of proposed system is presented in Fig. 1.

### III. EXPERIMENTAL SETUP

We evaluate the proposed SID module using the *speech separation challenge* corpus provided in [6]. The corpus is composed of 34 speakers (18 male, 16 female), with a total number of 34,000 utterances, each following a command-like structure, and all having a unique grammatical structure. Each sentence is formed by different syntaxes of command, color, letter, number and code, for instance "bin white by A 3 please". The test data in the corpus is composed of 500 laboratory-quality signals for each of the 34 target speakers, as well as test set consisting of mixed signals at six signal-to-signal ratio levels of  $\{-9, -6, -3, 0, 3, 6\}$  dB. For each of these six test sets for two-talker signal, 600 utterances are provided, from which 221 are for same talker (ST), 200 for same gender (SG), and 179 for different gender (DG). The utterances were originally sampled at 25 kHz with a duration of 2 second.

Since we are interested in telephone-quality speech bandwidth, we downsample the signals from 25 kHz to 8 kHz. We extract features from 30 msec frames multiplied by a Hamming window. A 27-channel mel-frequency filterbank is applied on discrete Fourier transform (DFT) spectrum to extract 12-dimensional mel-frequency cepstral coefficients (MFCCs), followed by appending  $\Delta$  and  $\Delta^2$  coefficients, and using an energy-based voice activity detector (VAD) for extracting the feature vectors. We digitally add the signals with an average frame-level SSR to construct the UBM

and the target speakers GMMs. For each of 34 speakers, 50 random files from each speaker were mixed at SSRs levels  $\{-9, -6, -3, 0, 3, 6\}$  dB with 50 random files from other speakers which gives us about 180 hour of speech for training UBM. The number of Gaussians,  $M$ , is set to 2048.

Speakers SSR-dependent GMMs,  $\lambda_{ig}$ , trained by mixing 100 random files from each speaker with 100 random files from other speakers yielding about 1.8 hours data for each SSR. Relevance factor was set to 16 for training speaker models,  $\lambda_{ig}$ , where its value was set to 0 in training test model,  $\lambda_e$ , because of availability of only 2 seconds of data for adaptation. We set the threshold  $\alpha$  to 1 in frame-level scores calculation. The accuracies defined here are to identify both of the speakers existing in mixed signal as the two most probable speakers.

### IV. EXPERIMENTAL RESULTS

We first analyze the performance of speaker identification system using each of the 5 scores individually. The results shown in Table I indicate that NWF and KLD have the best average performance compared to the other methods. To the best of our knowledge, SID accuracy for *Iroquois* is not reported without SSR estimation included. Compared to *LL* score, our proposed method, *NWF*, is more accurate. It is observed that, the number of frames above the confidence level, *NCF* is more important than their mean value, *LL*. On the other hand, the model based approach, *ACE*, works equally well as the frame-level method but it is more complex and has slightly worse accuracy than *KLD*.

Score fusion was then done by using two most successful methods:  $FUS_i = 0.54NWF_i + 0.46KLD_i$ . The fusion weights were optimized on development set consisting of 300 mixed signals for each SSR level. We found that, for the fusion system, in all of the experiments, one of the speakers in the mixed signal is *always* identified. The accuracy of the proposed system (*FUS*) for listing two target speakers in 3-best list is shown in Table II. This accuracy suggests to use proposed SID module as a concise "short-list" generator for the SSR estimation in *Iroquois* to reduce

Table I  
SPEAKER IDENTIFICATION ACCURACY FOR DIFFERENT SYSTEMS (PERCENTAGE OF UTTERANCES WITH BOTH SPEAKERS IN THE 2-BEST LIST OUTPUT). FUS IS PROPOSED SYSTEM COMPOSED OF 0.54NWF + 0.46KLD AND IRO STANDS FOR IROQUOIS

SSR (dB)	-9	-6	-3	0	3	6	Ave
NWF	81	90	94	95	92	88	90
NCF	75	88	93	94	92	86	89
LL	74	84	90	91	87	82	85
KLD	79	89	92	93	91	87	88
ACE	79	87	92	92	89	84	87
FUS	92	93	96	97	93	87	93
IRO [4]	96	98	98	99	99	98	98

Table II  
SPEAKER IDENTIFICATION ACCURACY FOR PROPOSED FUS SYSTEM (PERCENTAGE OF UTTERANCES WITH BOTH SPEAKERS IN THE 3-BEST LIST OUTPUT) ST, SAME TALKER, SG, SAME GENDER AND DG, DIFFERENT GENDER).

SSR	ST	SG	DG	Ave
-9 dB	100	93	83	92
-6 dB	100	97	94	97
-3 dB	100	100	98	99
0 dB	100	98	99	99
3 dB	100	97	93	97
6 dB	100	94	91	95
Ave	100	97	93	97

complexity. To understand the system performance better, we look for combinations of speakers that are identified in any given SSR. Surprisingly, in 68% of cases both speakers are correctly identified in the mixed signal at all SSR levels, and in 80% of experiments possibly only for one SSR we cannot identify both speakers but one of them. From the results, it is observed that mixed signals with different genders (DG) are more problematic than the same gender, which there are almost no significant difference in identification accuracy between males and females.

## V. CONCLUSION

A new method for speaker identification in co-channel scenario was introduced based on the existing approaches in speaker verification and compared the accuracy to *Iroquois* approach. From the simulation results conducted on speech separation challenge database, we observed that the proposed simple SID module performs well in listing two target speakers as three most probable speakers without any requirement on the estimates of the SSR level. As a future work, since we already got satisfactory results on 8 KHz speech, we plan to examine the proposed algorithm on telephony quality spontaneous speech and more realistically when signals are not synthetically mixed.

## REFERENCES

- [1] T. Kinnunen and H. Li, "An overview of text-independent speaker recognition: from features to supervectors," *Elsevier speech communication*, vol. 52, no. 1, pp. 12–40, January.
- [2] D. A. van Leeuwen, A. F. Martin, M. A. Przybicki, and J. S. Bouten, "NIST and NFI-TNO evaluations of automatic speaker recognition," *Elsevier Comp. Speech and Lang.*, vol. 20, no. 3, pp. 128–158, 2006.
- [3] D. Morgan, E. George, L. Lee, and S. Kay, "Co-channel speaker separation," in *Proc. IEEE Int. Conf. Acoust., Speech, Signal Processing*, 1995, pp. 828–831.
- [4] J. Hershey, S. Rennie, P. Olsen, and T. Kristjansson, "Super-human multi-talker speech recognition: A graphical modeling approach," *Elsevier Comp. Speech and Lang.*, vol. 24, no. 1, pp. 45–66, Jan 2010.
- [5] P. Mowlaee, R. Saeidi, Z. H. Tan, M. G. Christensen, P. Fränti, and S. H. Jensen, "Joint single-channel speech separation and speaker identification," in *Proc. IEEE Int. Conf. Acoust., Speech, Signal Processing*, March 2010, pp. 4430–4433.
- [6] M. Cooke, J. Hershey, and S. Rennie, "Monaural speech separation and recognition challenge," *Elsevier Comp. Speech and Lang.*, vol. 24, no. 1, pp. 1–15, 2010.
- [7] R. Weiss and D. Ellis, "Speech separation using speaker-adapted eigenvoice speech models," *Elsevier Comp. Speech and Lang.*, vol. 24, no. 1, pp. 16–29, 2010.
- [8] P. Li, Y. Guan, S. Wang, B. Xu, and W. Liu, "Monaural speech separation based on MAXVQ and CASA for robust speech recognition," *Elsevier Comp. Speech and Lang.*, vol. 24, no. 1, pp. 30–44, Jan. 2010.
- [9] D. Reynolds, T. Quatieri, and R. Dunn, "Speaker verification using adapted gaussian mixture models," *Digital Signal Processing*, vol. 10, no. 1-3, pp. 19–41, Jan 2000.
- [10] W. Campbell, D. Sturim, and D. Reynolds, "Support vector machines using GMM supervectors for speaker verification," *Signal Processing Letters, IEEE*, vol. 13, no. 5, pp. 308–311, May 2006.
- [11] H. Aronowitz and D. Burshtein, "Efficient speaker recognition using approximated cross entropy (ACE)," *Audio, Speech, and Language Processing, IEEE Transactions on*, vol. 15, no. 7, pp. 2033–2043, Sept. 2007.