



**HAL**  
open science

## Multi-stage speaker diarization of broadcast news

Claude Barras, Xuan Zhu, Sylvain Meignier, Jean-Luc Gauvain

► **To cite this version:**

Claude Barras, Xuan Zhu, Sylvain Meignier, Jean-Luc Gauvain. Multi-stage speaker diarization of broadcast news. *IEEE Transactions on Audio, Speech and Language Processing*, 2006, 14 (5), pp.1505-1512. 10.1109/TASL.2006.878261 . hal-01434241

**HAL Id: hal-01434241**

**<https://hal.science/hal-01434241v1>**

Submitted on 22 Mar 2017

**HAL** is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

# MULTI-STAGE SPEAKER DIARIZATION OF BROADCAST NEWS

Claude Barras\*, Xuan Zhu, Sylvain Meignier and Jean-Luc Gauvain, *Member, IEEE*

**Abstract**—This paper describes recent advances in speaker diarization with a multi-stage segmentation and clustering system, which incorporates a speaker identification step. This system builds upon the baseline audio partitioner used in the LIMSI broadcast news transcription system. The baseline partitioner provides a high cluster purity, but has a tendency to split data from speakers with a large quantity of data into several segment clusters. Several improvements to the baseline system have been made. First, the iterative Gaussian mixture model (GMM) clustering has been replaced by a Bayesian information criterion (BIC) agglomerative clustering. Second an additional clustering stage has been added, using a GMM-based speaker identification method. Finally a post-processing stage refines the segment boundaries using the output of a transcription system. On the NIST RT-04F and ESTER evaluation data, the multi-stage system reduces the speaker error by over 70% relative to the baseline system, and gives between 40% and 50% reduction relative to a single-stage BIC clustering system.

This work was partially financed by the European Commission under the FP6 Integrated Project IP 506909 CHIL and by the Defense Advanced Research Projects Agency under the GALE Program. Any opinions expressed in this paper are those of the authors and do not necessarily reflect the views of DARPA.

Claude Barras, Xuan Zhu and Jean-Luc Gauvain are with the LIMSI-CNRS, BP 133, 91403 Orsay, France. Tel: +33 1 6985 8061. Fax: +33 1 6985 8088. Mail: {barras,xuan.gauvain}@limsi.fr

Sylvain Meignier was with the LIMSI-CNRS and is now with the Laboratoire d'Informatique de l'Universite du Maine, France. Tel: +33 2 4383 3830. Fax: +33 2 4383 3868. Mail: sylvain.meignier@lium.univ-lemans.fr

**EDICS: SPE-SPKR Speaker Characterization and Recognition**

**Index Terms**—Speaker diarization, speaker identification, speaker segmentation and clustering, BIC clustering.

## I. INTRODUCTION

Speaker diarization, also called speaker segmentation and clustering, is the process of partitioning an input audio stream into homogeneous segments according to speaker identity. It is one aspect of audio diarization, along with categorization of music, background noise and channel conditions. Speaker diarization can improve the readability of an automatic transcription by structuring the audio stream into speaker turns and in some cases by providing the true speaker identity. Such information can also be of interest for the indexing of multimedia documents. As defined by NIST for the 2004 Rich Transcription evaluation [1], the speaker diarization task is relative to a given show and no a priori knowledge of the speaker's voice or even of the number of speakers is available. Therefore only a relative, show-internal speaker identification is produced by the diarization system. This definition has been adopted in this work, even though a speaker diarization system could obviously make use of such information if combined with a speaker identification or tracking system.

Audio diarization is a useful preprocessing step for an automatic speech transcription system. By separating out speech and non-speech segments, the recognizer only needs to process audio segments containing speech, thus reducing the computation time and avoiding word insertions in these portions. By clustering segments of the same acoustic nature, condition specific models can be used to improve the quality of the transcription. By clustering segments from the same speaker, the amount of data available for unsupervised speaker adaptation is increased, which can significantly improve the transcription performance.

Automatic speech transcription and speaker diarization rely on similar methods for segmentation and clustering. However differences in their objectives leads to different needs, particularly concerning where accuracy is most important. Automatic transcription requires accurate segment boundaries. Although the rejection of non-speech segments is useful in order to minimize insertion of words and to save computation time, it is important that the segment boundaries are located in non-informative zones such as silences or breaths. Indeed, having a word cut by a boundary disturbs the transcription process and increases the word error rate.

Diarization also aims to produce homogeneous speech segments; however, the main objectives are the purity and the correct labeling of the segments. Errors such as having more than one cluster for a given speaker, or conversely, merging the segments of two different speakers into one cluster, are penalized more heavily. The effects of both boundary inaccuracy and mislabeled segments were measured in [2] for English broadcast news. The experiments showed that segment boundary errors have a greater impact on the transcription task, while label errors have a greater impact on the diarization task.

The NIST Rich Transcription evaluation has been the major evaluation for speaker diarization of broadcast news data in 2003 and 2004 [1], [3], [4]. In 2005, the TechnolanguagE ESTER evaluation was conducted on a similar task using French radio broadcast news data [5], [6].

Most speaker diarization systems for BN data have a similar general architecture. First the signal is chopped into homogeneous segments. The segment boundaries are located by finding acoustic changes in the signal and each segment is expected to contain speech from only one speaker. The resulting segments are then clustered so that each cluster corresponds to one speaker, a major issue being that the number of speakers is unknown a priori and needs to be automatically determined. Each system also presents specific aspects which can be classified following different criteria:

- Link between segmentation and clustering: segmentation can be done first, followed by clustering with no connection between the two parts [7], [8] inspired from the work presented in [9]–[11]; alternatively the segmentation and clustering can be jointly optimized, via, for example, the iterative segmentation and clustering procedures described in [12]–[14]. A limitation of the first method is that errors made in the segmentation step are not only difficult to correct later, but can also degrade the performance of the subsequent clustering step.
- Clustering strategy: it relies either on an agglomerative clustering [7], [12], [14], or on a divisive clustering method [13], [15].
- Modeling strategy: each speaker can be modeled by a Gaussian Mixture Model (GMM) with diagonal covariance matrices composed of 8 to 64 components. As is done in the speaker recognition task,

larger models with 2048 components have been proposed [8], [13], [14]. In this case, a more robust estimation of the models despite the limited amount of data per speaker can be obtained by performing the maximum a posteriori (MAP) adaptation of a prior model [16]. On the other hand, using a single Gaussian with a full covariance matrix for the modeling of a speaker also provides good results [7].

In our experiments several variants and combinations of systems have been tested, in particular to study the link between segmentation and clustering and the modeling strategy.

The remainder of this paper is organized as follows: Section II describes the baseline partitioning system which was developed for the automatic Broadcast News transcription task. Section III describes the multi-stage partitioning system specifically built for the speaker diarization task. This system is based upon a Bayesian information criterion (BIC) clustering followed by speaker identification (SID) clustering. Experimental results are presented in Section IV, followed by some conclusions.

## II. BASELINE PARTITIONING SYSTEM

The baseline audio partitioning system was developed as a preprocessing step for the LIMSI English broadcast news transcription system [12], [17]. It was shown to provide a high cluster purity (about 96%) and a cluster coverage slightly below 80% on 1996 and 1997 NIST evaluation data. This baseline partitioner **c-std** shown in Figure 1 is structured as follows:

### A. Feature extraction

Mel frequency cepstral parameters are extracted from the speech signal every 10 ms using a 30 ms window on a 0-8kHz band. For each frame the Mel scale power spectrum is computed, and the cubic root taken followed

by an inverse Fourier transform. Using a process similar to that of PLP computation [19], 12 LPC-based cepstral coefficients are then extracted. The 38 dimensional feature vector consists of 12 cepstral coefficients,  $\Delta$  and  $\Delta$ - $\Delta$  coefficients plus the  $\Delta$  and  $\Delta$ - $\Delta$  log-energy. This is essentially the same set of features that is used in a standard transcription system, except for the energy [18]. This set is used in all steps of the **c-std** system, except for the segmentation into small segments where only the static features are used. No cepstral mean or variance normalization is performed to the acoustic vector in the baseline partitioning system.

### B. Speech Activity Detection (SAD)

Speech is extracted from the signal with a Viterbi decoding using Gaussian Mixture Models (GMM) for speech, noisy speech, speech over music, pure music, and silence or noise. The aim of the SAD is to remove only long regions without speech such as silence, music and noise, so the penalty of switching between models in the Viterbi decoding was set to minimize the loss of speech signal. The GMMs, each with 64 Gaussians, were trained on about 1 hour of the specific type of data, selected from English Broadcast News data from 1996 and 1997 distributed by the LDC (Linguistic Data Consortium).

### C. Chopping into small segments

The segmentation process consists of finding segment boundaries that correspond to the instantaneous speaker change points. It is generally a good choice to minimize the miss rate for speaker change points even if the false alarm rate is high, because the false change points can be easily removed later during a clustering process. The segmentation needs to provide pure segments (i.e.

containing speech from only one speaker) of duration sufficient to characterize the voice of the speaker.

Segmentation of the signal is performed by taking the maxima of a local Gaussian divergence measure between two adjacent sliding windows  $s_1$  and  $s_2$ . For each segment, the static features (i.e., only the 12 cepstral coefficients plus the energy) are modeled with a single diagonal Gaussian, i.e.  $s_1 \sim \mathcal{N}(\mu_1, \Sigma_1)$  and  $s_2 \sim \mathcal{N}(\mu_2, \Sigma_2)$  with  $\Sigma_1$  and  $\Sigma_2$  diagonal. Then the Gaussian divergence measure is defined as:

$$G(s_1, s_2) = (\mu_2 - \mu_1)' \Sigma_1^{-1/2} \Sigma_2^{-1/2} (\mu_2 - \mu_1) \quad (1)$$

It is the Mahalanobis distance between  $\mu_1$  and  $\mu_2$  weighted by the geometric mean of  $\Sigma_1$  and  $\Sigma_2$ , which reduces to a weighted Euclidean distance because of the diagonal assumption. The detection threshold was optimized on the training data in order to provide acoustically homogeneous segments. The window size was set to 5 seconds with a minimal segment length of 2.5 seconds. Due to the simple diagonal assumption, this segmentation phase is very quick. This approach is similar to the segmentation proposed in [9] using the symmetric KL2 metric. Other popular segmentation methods are based upon the BIC metric [10], [20], [21], but these methods show a much higher complexity [22]. An analysis of various speaker change point techniques based on models, metrics or energy is given in [23].

#### D. Iterative GMM segmentation/clustering procedure

Each initial segment is used to seed one cluster, and an 8-component GMM with a diagonal covariance matrix is trained on the segment's data. The algorithm alternates the Viterbi resegmentation and the GMM re-estimation and merging steps, with the goal of maximizing the objective function:

$$\sum_{i=1}^N \log f(s_i | M_{c_i}) - \alpha N - \beta K \quad (2)$$

where  $S = (s_1, \dots, s_N)$  is the partitioning of the speech segments into a sequence of  $N$  segments,  $c_i \in [1, K]$  is the cluster label for the segment  $s_i$  among the  $K$  different clusters,  $f(s_i | M_{c_i})$  is the likelihood of the segment  $s_i$  given the model of its cluster  $M_{c_i}$ , and  $\alpha$  and  $\beta$  are the segment and cluster penalties. The procedure stops when no more merges are possible. More details on the clustering procedure can be found in [12]. This procedure is similar to BIC using a global penalty as described in Section III-A.

#### E. Viterbi resegmentation

The segment boundaries are refined using the last set of GMMs and an additional relative energy-based boundary penalty, within a one second interval. The boundaries are thus shifted to the nearest point of low energy within this interval. This is done so as to locate the segment boundaries at silence portions, thereby avoiding cutting words. This is especially important when using the resulting segmentation as a pre-processing step of an automatic transcription system.

#### F. Bandwidth and gender labeling

Bandwidth (wide band studio or narrow band telephone) detection for each segment is first performed using two GMMs. The gender (male or female) labeling is then carried out on the segments using two pairs of bandwidth dependent GMMs: for gender labeling on telephone speech segments, feature extraction is limited to the 0-3.5kHz band. The GMM models are composed of 64 components with diagonal covariance matrices and were trained on the subset of the LDC 1996/1997 English Broadcast News data also used to train the

speech detection models. This labeling is useful for the transcription system, as different acoustic models are used for each combination of bandwidth and gender for better performance, but is also of interest for structuring the acoustic stream. Performing the labeling on a segment basis rather than for a whole cluster may split a cluster in two which can prove beneficial, since a given speaker is usually not recorded in both wide band and narrow band conditions in the same show.

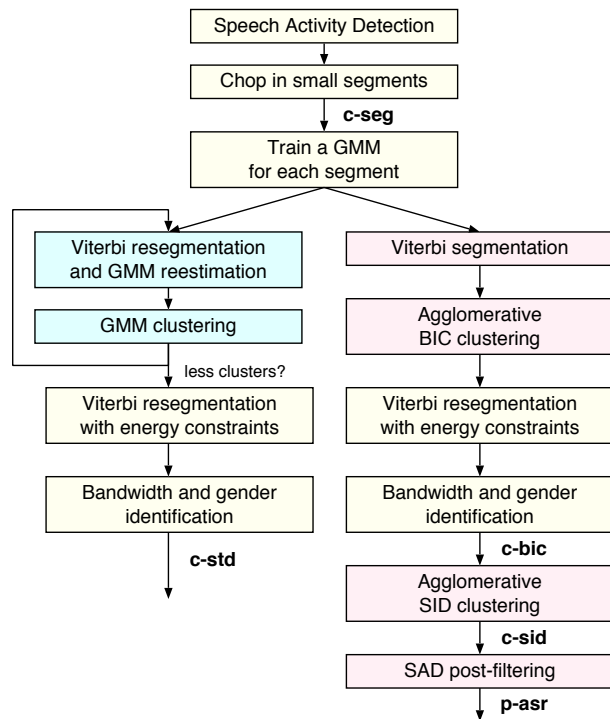


Fig. 1. Architecture of the baseline partitioning system (**c-std** on the left side of the diagram) and the multi-stage speaker diarization system (**p-asr** to the right, along with **c-seg**, **c-bic** and **c-sid** intermediate steps).

### III. MULTI-STAGE DIARIZATION

Recent research has shown BIC clustering methods to obtain good performance on the speaker diarization

task [7], [15]. We therefore tested a modified system, replacing the iterative GMM clustering with a BIC-based clustering (cf. Figure 1) and using Gaussian models with a full covariance matrix. We believe the iterative resegmentation to be more important in the context of speaker partitioning for the transcription task than for the speaker diarization task. Since different models can capture different and complementary aspects of the data, we decided to combine them in a multi-stage system. We decided to pipeline the output of the system with the BIC clustering into a second clustering stage which uses a speaker identification module. The SID clustering uses a more aggressive acoustic channel normalization and a more complex speaker model, enabled by the larger amount of data per cluster at the beginning of this stage. Finally, an SAD post-filtering stage was added in order to remove short intra-speaker pauses. These short pauses while indeed harmless for a speech transcription system, are penalized as false alarms in a speaker diarization system. The other parts of the system were kept unchanged.

#### A. BIC clustering

Agglomerative clustering is applied to the segments resulting from the GMM segmentation. Initially, each segment seeds one cluster, modeled by a single Gaussian with a full covariance matrix trained on the 12 Mel frequency cepstrum coefficients and the energy (but without the  $\Delta$  coefficients). At each iteration, the two nearest clusters are merged until the stopping criterion is reached. The BIC criterion [10] is used both for the inter-cluster distance measure and the stop criterion.

In order to decide whether to merge two clusters  $c_i$  and  $c_j$ , the  $\Delta BIC$  value is computed as:

$$\Delta BIC = (n_i + n_j) \log |\Sigma| - n_i \log |\Sigma_i| - n_j \log |\Sigma_j| - \lambda P \quad (3)$$

where  $\Sigma$  is the covariance matrix of the merged cluster ( $c_i$  and  $c_j$ ),  $\Sigma_i$  of cluster  $c_i$ ,  $\Sigma_j$  of cluster  $c_j$ , and  $n_i$  and  $n_j$  are respectively the number of the acoustic frames in clusters  $c_i$  and  $c_j$ . The penalty  $P$  is:

$$P = \frac{1}{2}(d + \frac{1}{2}d(d + 1)) \log n \quad (4)$$

where  $d$  is the dimension of the feature vector space. The term  $n_i \log |\Sigma_i|$  is related to the log likelihood of the cluster  $c_i$  given its estimated Gaussian  $M_{c_i}$ <sup>1</sup>. Singular covariance matrices were not an issue because of the minimal length constraint during the initial segmentation. The merging criterion is that two clusters should be merged if  $\Delta BIC < 0$ . At each step the two nearest clusters (i.e. those which have the most negative  $\Delta BIC$  value) are merged into one cluster, and the  $\Delta BIC$  values between this new cluster and remaining clusters are computed. This clustering procedure terminates when the  $\Delta BIC$  between all cluster pairs is greater than zero.

In our BIC clustering procedure, the size of the two merged clusters, i.e.  $n = n_i + n_j$ , is used to compute the penalty  $P$ , as described in [21]. We refer to this as a local BIC penalty. Another solution is to use the size of the whole set of clusters, i.e.  $n = \sum_{k=1}^N n_k$  to compute the penalty, which we refer to as a global BIC penalty and corresponds to an exponential prior for the number of clusters. In this case the penalty is constant and the decision to merge two clusters is decided just by the increase in likelihood. This in fact corresponds to the objective function in Equation (2) used in the baseline partitioner when the number of segments is fixed. For broadcast news documents, our experimental results as presented in Section IV demonstrate the local BIC to be a better choice for a merging criterion.

<sup>1</sup>more precisely,  $\log f(c_i|M_{c_i}) = -\frac{n_i}{2} \log |\Sigma_i| - \frac{n_i d}{2}(1 + \log 2\pi)$ , but the constant factor  $\frac{1}{2}$  was simplified in Equation 3.

## B. SID clustering

After the initial segmentation, both the iterative GMM and the agglomerative BIC clustering methods have to deal in the beginning of the process with short duration segments, and thus use a limited set of parameters per cluster. After several iterations, the amount of data per cluster increases, so a more complex model can be used. Our approach is to stop the initial clustering stage early, and use the results to seed a second clustering stage with more initial data per cluster. This second stage can therefore estimate more complex models for the speakers. In addition, purely acoustic clustering tends to split a speaker's data into several clusters as a function of the various background conditions (clean speech, speech with noise, speech with music etc.), so an acoustic background normalization is necessary to regroup the data for a given speaker.

After the BIC clustering stage, state-of-the-art speaker recognition methods [24], [25] were used to improve the quality of the speaker clustering. The feature vector consists of 15 Mel frequency cepstral coefficients plus delta coefficients and delta energy. Feature warping normalization, which reshapes the histogram of the cepstral coefficients into a Gaussian distribution [26] is performed on each segment using a sliding window of 3 seconds in order to reduce the effect of the acoustic environment.

For each gender and channel condition (studio, telephone) combination, a Universal Background Model (UBM [27]) with 128 diagonal Gaussians is trained on the 1996/1997 English Broadcast News data. The GMM for each remaining cluster is obtained by maximum a posteriori (MAP) adaptation [16] of the means of the matching UBM.

Agglomerative clustering is performed separately for

each gender and bandwidth condition, using a cross log-likelihood ratio as in [28]. For each cluster  $c_i$ , its model  $M_i$  is MAP adapted from the gender and channel matched UBM  $B$  using the feature vectors  $x_i$  belonging to the cluster. Given two clusters  $c_i$  and  $c_j$ , the cross log-likelihood ratio  $\mathcal{S}$  is defined as:

$$\mathcal{S}(c_i, c_j) = \frac{1}{n_i} \log \frac{f(x_i|M_j)}{f(x_i|B)} + \frac{1}{n_j} \log \frac{f(x_j|M_i)}{f(x_j|B)} \quad (5)$$

where  $f(\cdot|M)$  is the likelihood of the acoustic frames given the model  $M$ , and  $n_i$  is the number of frames in cluster  $c_i$ . This is a symmetric similarity measure. After each merge, a new model is trained for the cluster  $c_i \cup c_j$ . The clustering stops when the cross log-likelihood ratio between all clusters is below a given threshold  $\delta$  optimized on the development data.

### C. SAD post-filtering

In order to filter out short-duration silence segments that were not removed in the initial speech detection step to further reduce the speaker diarization error, a post-processing stage uses the word segmentation output by the LIMSI Broadcast News Speech-To-Text system [29] relying on the **c-std** system for the segmentation and clustering. Only inter-word silences shorter than 1 second are filtered out, this value being determined empirically.

## IV. EXPERIMENTS AND RESULTS

Experiments are reported for diarization systems that were submitted to the NIST and ESTER evaluations [30]. Several configurations were tested for the systems. Unless otherwise specified, the configuration used is the one that provided the best results on development data, i.e.  $\alpha = \beta = 230$  for **c-std**,  $\lambda = 5.5$  for **c-bic** and  $\lambda = 3.5, \delta = 0.1$  for **c-sid** and **p-asr**. In the **c-sid** system, the BIC penalty weight  $\lambda$  was optimized to

cluster only the closest segments in the BIC clustering stage so as to give more degrees of freedom to the SID clustering stage; while in the **c-bic** system,  $\lambda$  was optimized directly to give the lowest diarization error the BIC clustering could bring. A local BIC merging and stop criterion was also used.

### A. Corpora

The experiments were conducted on the US English data used in NIST RT-04F (Fall 2004 Rich Transcription Evaluation) [1] and on the French data from the French ESTER broadcast news evaluation [5].

For the RT-04F evaluation, the training and development corpora were provided for system development along with a reference speaker labeling determined by LDC. Evaluation references were made available after the evaluation. The data are from US radio or TV broadcast news shows. The development data has two portions, '*dev1*' and '*dev2*' each consisting of 6 30-minute audio files. *Dev1* was recorded in February 2001, with programs from ABC, CNN, NBC, PRI and VOA; and *dev2* was recorded in November and December 2003, with programs from ABC, CNBC, CNN, CSPAN and PBS. The evaluation data is comprised of 12 audio files each lasting approximately 30 minutes, recorded in December 2003, extracted from the shows of ABC, CNBC, CNN, CSPAN, PBS and WBN. The *dev2* and evaluation and corpora are very similar to each other since the shows are recorded from the same channels, at the same period, whereas *dev1* is an older corpus coming from the previous RT-03 evaluation.

The data used in ESTER were extracted from French radio broadcast news shows, provided by ELDA (Evaluations and Language resources Distribution Agency) and the DGA (Délégation Générale pour l'Armement). The training corpus contains 82 hours of data from the France



Inter, France Info, RFI and RTM radio stations, recorded in 1998, 2000 and 2003, with the audio file durations ranging from 10 minutes to 1 hour. The development corpus contains 8 hours of data, in 14 audio files recorded from April to July 2003 from the same stations as the training corpus. The evaluation corpus is comprised of 18 audio files recorded from October to December 2004, with a total duration of 10 hours. The evaluation corpus contains data from two radio stations ('France Culture' and 'Radio Classique') not present in the training or development corpora. There is also a 14 month interval between the recording period of the development and test data (July 2003 to October 2004) of these two corpora. Both data sets present a large variability in audio file durations (10 minutes to 1 hour).

### B. Performance measures

The speaker diarization task performance is measured via an optimum one-to-one mapping between the reference speaker IDs and the hypothesis speaker IDs. The primary metric for the task, referred to as the speaker match error, is the fraction of speaker time that is not attributed to the correct speaker, given the optimum speaker mapping. Another measure is the overall speaker diarization error rate (DER) which includes the missed and false alarm speaker times, thus taking speech/non-speech detection errors into account [1]. All the measurements mentioned above are illustrated in Figure 2.

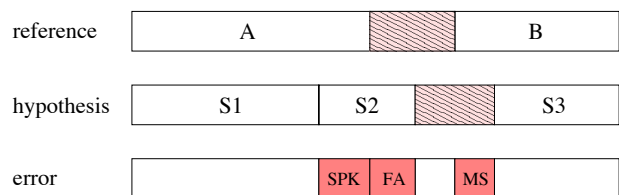
In order to more closely analyze the performance of speaker clustering methods, average frame-level cluster purity and cluster coverage are used as defined by [12]. Cluster purity is defined as the ratio between the number of frames by the dominating speaker in a cluster and the total number of frames in the cluster. Cluster coverage is a dual measure, and accounts for the dispersion of a given speaker's data across clusters; for a given speaker,

it is defined as the percentage of its frames in the cluster which has most of the data of the speaker. Cluster coverage was also expressed as the purity of reference clusters in [21]. In these experiments, cluster purity and cluster coverage errors are reported.

Cluster purity and coverage measures are complementary, and the speaker match error time can be interpreted as a combination of both. Moreover, it is interesting to note that if, the hypothesized speaker for a cluster is always taken to be the majority reference speaker in this cluster, then the speaker match error will be exactly the cluster purity error; it is easy to demonstrate that it is also a lower bound for the match error. Thus, starting with an initial segmentation, the cluster purity error will be the lowest possible match error on this segmentation after performing an agglomerative clustering. The same holds for further clustering of the output of a previous clustering stage, as long as the segment boundaries are not modified.

### C. Results on the RT-04F development data

The performance at different stages of the system was compared for different system configurations, as presented in Table I. On RT-04F *dev1* data, the initial segmentation **c-seg**, with the minimal duration constraint of 2.5 sec per segment, has a purity error of 1% which is also the lowest possible speaker match error. The



$$\text{DER} = \text{Speaker Error (SPK)} + \text{False Alarm Speech (FA)} + \text{Missed Speech (MS)}$$

Fig. 2. Example of the performance measures for the speaker diarization task used in the NIST RT-04F evaluation.

coverage error of this initial segmentation is of course very high at 73.2%; which means that on average, only about one quarter of the speech for each speaker is located in a single segment.

As expected, the standard partitioner **c-std** in its default configuration provides good purity, but relatively poor coverage, resulting in a high overall diarization error of over 30% on *dev1* data. Setting the penalty  $\alpha$  and  $\beta$  to optimize these values reduces this error below 25%. The **c-bic** system also provides a high purity, with a much better coverage (2.9% purity error and 9.8% coverage error), reducing the overall error rate by almost 50%. The **c-sid** system achieves a large decrease of the coverage error with a further small improvement of the purity, resulting in an overall DER of 7.1%, a reduction of almost 50% compared to the **c-bic** system.

A global BIC merging and stop criterion was also tested, but always performed worse than the local BIC criterion in our experiments, as can be seen for **c-bic** system on RT-04F *dev1* in Table II. A similar result was observed in [15]. This result needs further investigation, but may be due to a mismatch between the BIC model and the real distribution of the data. Thus only the local criterion was used in the remaining experiments.

The effect of the SID detection threshold  $\delta$  on the speaker match error and the cluster purity error was measured on both the *dev1* and *dev2* data. A lower threshold reduces the number of final speaker clusters. As shown in Figure 3, reducing the threshold in the positive range results in a decrease of the match error rate with almost no degradation of the purity. Moreover, there is a large range of thresholds around zero with a low speaker match error. The cluster purity error shows the speaker match error that could be achieved with the best clustering decision, as explained in section IV-B.

Looking at the performance of the **c-sid** system in

TABLE I  
THE PURITY, COVERAGE AND OVERALL DIARIZATION ERROR RATES FROM THE **c-std**, **c-bic** AND **c-sid** SYSTEMS ON THE RT-04F AND THE ESTER DEVELOPMENT DATASETS.

<i>system</i>	<i>purity error (%)</i>	<i>coverage error (%)</i>	<i>overall DER</i>
RT-04F <i>dev1</i> data set			
c-seg	1.0	73.2	N/A
c-std ( $\alpha = \beta = 160$ )	5.0	28.4	32.3
c-std ( $\alpha = \beta = 230$ )	9.4	17.9	24.8
c-bic ( $\lambda = 5.5$ )	2.9	9.8	13.2
c-sid ( $\lambda = 3.5, \delta = 0.1$ )	2.1	4.2	7.1
RT-04F <i>dev2</i> data set			
c-sid ( $\lambda = 3.5, \delta = 0.1$ )	1.7	3.5	7.6
ESTER development data set			
c-bic ( $\lambda = 5.5$ )	7.2	10.6	15.8
c-sid ( $\lambda = 3.5, \delta = 1.5$ )	4.7	5.2	8.0

TABLE II  
THE OVERALL DIARIZATION ERROR FOR **c-bic** SYSTEM ON THE RT-04F *dev1* DATA, AS A FUNCTION OF THE PENALTY WEIGHT  $\lambda$  FOR THE LOCAL AND GLOBAL BIC CRITERION.

<i>BIC criterion</i>	$\lambda$	<i>overall DER (%)</i>	<i>BIC criterion</i>	$\lambda$	<i>overall DER (%)</i>
local	5.0	13.3	global	5.0	16.4
	6.0	12.8		6.0	15.5
	7.0	13.8		7.0	18.2

more detail, a large variation in the speaker error is observed across shows, as shown in Table III. The speaker error ranges from a low of 0.1% to over 12%. Having only very few speakers (3), the CSPAN show has the lowest speaker error. The ABC and NBC shows have more speakers, occurring in different background conditions, which is more challenging for the diarization system.

#### D. Results on the ESTER training and development data

For the French ESTER data, SAD was performed using the same American English speech/non-speech acoustic models as were used for RT-04F plus an additional speech over music model trained on French broadcast news data for a better recognition of the jingles found in the French radio data. The optimal threshold for the SID clustering on development data was  $\delta = 1.5$ . As can be seen in Table I, the **c-sid** system also has low purity and coverage error rates (4.7% and 5.2%, respectively) on the ESTER development data. A 50% reduction of the overall error rate is gained by adding the **c-sid** system to the **c-bic** system. The  $\delta$  threshold found for the RT-04F data did not carry over to the ESTER data, and this may be due to the larger variability in show sources, durations and types observed in ESTER.

The training data was divided into four subsets according to the show duration (10, 15, 20 minutes, 1 hour). As shown in Table IV, different optimal values of SID clustering threshold  $\delta$  were obtained for each subset. Larger SID clustering thresholds are better for the longer shows. The optimal SID clustering threshold

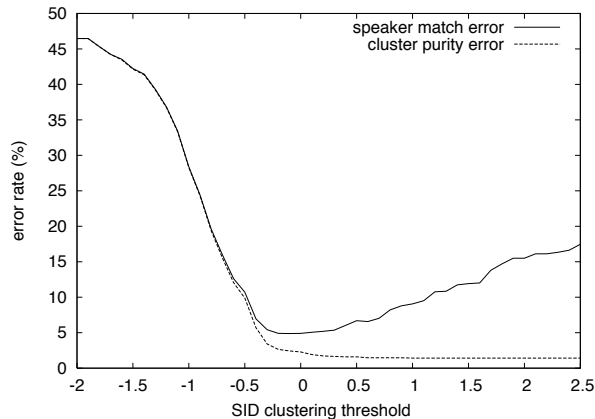


Fig. 3. Speaker match error and purity error rates on RT-04F *dev1* and *dev2* for the **c-sid** system as a function of the SID clustering threshold  $\delta$ .

TABLE III

PERFORMANCE OF **c-sid** SYSTEM ON THE RT-04F DEVELOPMENT DATASET, SCORES ARE GIVEN FOR MISS (MS), FALSE ALARM (FA), SPEAKER ERROR (SPK) AND OVERALL DIARIZATION ERROR RATE (DER), #REF AND #SYS ARE RESPECTIVELY THE REFERENCE AND SYSTEM SPEAKER NUMBER.

<i>show</i>	<i>#REF</i>	<i>#SYS</i>	<i>MS</i>	<i>FA</i>	<i>SPK</i>	<i>DER</i>
<b>dev1</b>	121	161	<b>0.4</b>	<b>1.3</b>	<b>5.4</b>	<b>7.1</b>
ABC	27	37	1.6	1.3	12.4	15.2
VOA	20	22	0.3	1.2	2.2	3.7
PRI	27	30	0.1	0.9	2.8	3.8
NBC	21	35	0.1	1.1	12.0	13.2
CNN	16	21	0.5	1.4	5.6	7.6
MNB	10	16	0.2	1.8	0.8	2.8
<b>dev2</b>	90	130	<b>0.5</b>	<b>3.1</b>	<b>4.1</b>	<b>7.6</b>
CSPAN	3	4	0.3	2.9	0.1	3.3
CNN	17	22	0.6	4.2	5.0	9.8
PBS	27	29	0.1	2.8	7.4	10.3
ABC	23	29	2.1	6.7	12.5	21.2
CNNHL	9	26	0.0	1.4	0.5	1.9
CNBC	11	20	0.2	1.0	0.9	2.1

$\delta$  on all of the training data is the same as the one for the development data. The show duration was taken here as a rough indicator of the speaker count, however a more appropriate model of the show type would enable a finer analysis.

#### E. Results on the evaluation data

The trends observed on the development data were confirmed on the RT-04F evaluation data. The results given in Table V, show a slight increase in overall diarization error to 17% for the **c-bic** system and to 9.1% for the **c-sid** system. The final SAD post-processing stage gives an improvement of 0.6%, mainly by reducing false alarms in speech detection. As mentioned in [31], the **p-asr** system had the best performance of all the

TABLE IV  
RESULTS OF THE DIFFERENT OPTIMAL SID CLUSTERING  
THRESHOLD  $\delta$  FOR THE ESTER TRAINING SUBSETS WITH THE  
DIFFERENT SHOW DURATIONS.

<i>training subset</i>	$\delta$	<i>purity error (%)</i>	<i>coverage error (%)</i>	<i>speaker match error (%)</i>
10 min	0.9	1.2	1.5	2.6
15 min	0.9	3.6	0.5	3.8
20 min	1.1	2.2	1.6	3.5
1 hour	1.5	3.5	5.2	7.6
<b>all</b>	<b>1.5</b>	<b>3.0</b>	<b>4.7</b>	<b>6.7</b>

TABLE V  
PERFORMANCES OF **c-bic**, **c-sid** AND **p-asr** SYSTEMS ON THE  
EVALUATION DATA OF RT-04F AND ESTER.

<i>system</i>	<i>missed speech</i>	<i>false alarm speech</i>	<i>speaker error</i>	<i>overall DER</i>
RT-04F test dataset				
c-bic	0.4	1.8	14.8	17.0
c-sid( $\delta = 0.1$ )	0.4	1.8	6.9	9.1
p-asr	0.6	1.1	6.8	8.5
ESTER test dataset				
c-bic	0.7	1.0	12.1	13.8
c-sid( $\delta = 1.5$ )	0.7	1.0	9.8	11.5
post-evaluation result on ESTER test dataset				
c-sid( $\delta = 2.0$ )	0.7	1.0	7.4	9.1

participants of the RT-04F evaluation by a significant margin.

Results on the ESTER evaluation data are given in Table V, with the setting optimized on the development data. The overall diarization error was reduced from 13.8% for the **c-bic** system to 11.5% for the **c-sid** system. The submitted system also had the best performance for this task in the ESTER evaluation [6]. In a post-evaluation experiment, a 20% relative reduction of the

overall diarization error was observed for the **c-sid** system with the best a posteriori threshold, showing that the error rate is highly dependent on the clustering threshold.

## V. CONCLUSIONS

In this paper, a multi-stage architecture for speaker diarization of broadcast news has been described. It builds upon a baseline speaker partitioning system which had been optimized for the automatic transcription task, but the constraints of the speaker diarization task are different. Several modifications to the baseline system have thus been explored. First, the iterative GMM clustering was replaced by an agglomerative BIC clustering, using single full-covariance Gaussian models. A local BIC merging and stop criterion was shown to outperform the global criterion; a similar result was observed in [15]. A second clustering module was applied to the output of the system, relying on techniques used for speaker identification and verification: acoustic channel normalization, and MAP adaptation of a reference GMM with a large number of Gaussians.

The multi-stage system was demonstrated to perform much better than the baseline system for the diarization task. On the RT-04F development data, a relative error reduction of over 70% was achieved when compared to the baseline system. This system obtained the best diarization performance in both the NIST RT-04F and the ESTER evaluations by a significant margin. An overall diarization error rate under 10% was obtained on the RT-04F evaluation data, while the error rate of the BIC-based systems were over 15%. Focusing on the speaker error only, the multi-stage system provides an error reduction of up to 50% relative to a standard BIC clustering system. This dramatic improvement over the baseline system results from several changes: the combination of

two different clustering stages, each one focusing on a different acoustic aspect with more complex modeling in the second stage, and the use of acoustic channel normalization methods suited to speaker identification. A system following this architecture recently developed at Cambridge University demonstrated similar improvements [32], where it was observed that a very important part of the gain was obtained by the feature warping normalization.

Several remaining issues need further investigation in order to improve the robustness and the efficiency of the system. It was observed that the clustering threshold needs to be set according to the type of the audio document, and that the system still has a large variability across individual shows. Only with a large amount of files can statistically consistent results be obtained. This is especially important since the speaker error does not provide a stable and continuous measure of a clustering system: splitting a speaker in two classes, which is a single decision, results in doubling of the error rate for this speaker.

Finally, most speaker diarization systems rely on a purely acoustic segmentation and clustering, whereas an essential part of the information in speech is of a linguistic nature, and obviously in TV and radio shows most speakers are presented and identified. Combining the acoustic information with the linguistic layer as explored in [33] would improve the robustness of a speaker diarization system and make it more exploitable by a human reader.

#### REFERENCES

- [1] NIST, "Fall 2004 Rich Transcription (RT-04F) evaluation plan," <http://www.nist.gov/speech/tests/rt/rt2004/fall/docs/rt04f-eval-plan-v14.pdf>, August 2004.
- [2] S. E. Tranter, K. Yu, D. A. Reynolds, G. Evermann, D. Y. Kim, and P. C. Woodland, "An investigation into the interactions between speaker diarisation systems and automatic speech transcription," Cambridge University Engineering Department Technical Report CUED/F-INFENG/TR-464, Tech. Rep., October 2003.
- [3] NIST, "The Rich Transcription Spring 2003 (RT-03S) evaluation plan," <http://www.nist.gov/speech/tests/rt/rt2003/spring/docs/rt03-spring-eval-plan-v4.pdf>, February 2003.
- [4] —, "Spring 2004 (RT-04S) Rich Transcription meeting recognition evaluation plan," <http://www.nist.gov/speech/tests/rt/rt2004/spring/documents/rt04s-meeting-eval-plan-v1.pdf>, February 2004.
- [5] G. Gravier, J.-F. Bonastre, S. Galliano, E. Geoffrois, K. Mc Tait, and K. Choukri, "The ESTER evaluation campaign of rich transcription of French broadcast news," in *Language Evaluation and Resources Conference (LREC 2004)*, Lisbon, Portugal, May 2004.
- [6] S. Galliano, E. Geoffrois, D. Mostefa, K. Choukri, J.-F. Bonastre, and G. Gravier, "The ESTER phase II evaluation campaign for the rich transcription of French broadcast news," in *Proceedings of the 9th European Conference on Speech Communication and Technology (ISCA Interspeech'05)*, Lisbon, Portugal, September 2005, pp. 1149–1152.
- [7] Y. Moh, P. Nguyen, and J.-C. Junqua, "Towards domain independent speaker clustering," in *Proceedings of International Conference on Acoustics Speech and Signal Processing (IEEE, ICASSP 2003)*, Hong Kong, April 2003.
- [8] M. Ben, M. Betsler, F. Bimbot, and G. Gravier, "Speaker diarization using bottom-up clustering based on a parameter-derived distance between GMMs," in *Proceedings of International Conference on Spoken Language Processing (ISCA, ICSLP 2004)*, Jeju, Korea, October 2004.
- [9] M. Sieglar, U. Jain, B. Raj, and R. Stern, "Automatic segmentation and clustering of broadcast news audio," in *the DARPA Speech Recognition Workshop*, Chantilly, VA, USA, February 1997.
- [10] S. Chen and P. Gopalakrishnan, "Speaker, environment and channel change detection and clustering via the Bayesian information criterion," in *DARPA Broadcast News Transcription and Understanding Workshop*, Landsdowne, VA, USA, February 1998.
- [11] T. Hain and P. C. Woodland, "Segmentation and classification of broadcast news audio," in *Proceedings of International Conference on Spoken Language Processing (ICSLP 98)*, Sydney, Australia, November 1998.
- [12] J.-L. Gauvain, L. Lamel, and G. Adda, "Partitioning and transcription of broadcast news data," in *Proceedings of International*

- Conference on Spoken Language Processing (ICSLP 98)*, Sydney, Dec. 1998, pp. 1335–1338.
- [13] S. Meignier, J.-F. Bonastre, and S. Igounet, “E-HMM approach for learning and adapting sound models for speaker indexing,” in *2001: a Speaker Odyssey. The Speaker Recognition Workshop (ISCA, Odyssey 2001)*, Chania, Crete, June 2001, pp. 175–180.
- [14] J. Ajmera and C. Wooters, “A robust speaker clustering algorithm,” in *Automatic Speech Recognition and Understanding (IEEE, ASRU 2003)*, St. Thomas, U.S. Virgin Islands, November 2003, pp. 411–416.
- [15] S. E. Tranter and D. A. Reynolds, “Speaker diarisation for broadcast news,” in *2004: A Speaker Odyssey. The Speaker Recognition Workshop (ISCA, Odyssey 2004)*, Toledo, Spain, May 2004.
- [16] J.-L. Gauvain and C. H. Lee, “Maximum a posteriori estimation for multivariate Gaussian mixture observations of Markov chains,” *IEEE Transactions on Speech and Audio Processing*, vol. 2(2), pp. 291–298, April 1994.
- [17] J.-L. Gauvain, L. Lamel, and G. Adda, “Audio partitioning and transcription for broadcast data indexation,” *Multimedia Tools and Applications*, vol. 14, pp. 187–200, 2001.
- [18] —, “The LIMSI broadcast news transcription system,” *Speech Communication*, vol. 37, no. 1-2, pp. 89–108, 2002.
- [19] H. Hermansky, “Perceptual linear predictive (PLP) analysis of speech,” *J. Acoustic. Soc. America*, vol. 87(4), pp. 1738–1752, 1990.
- [20] P. Delacourt and C. Wellekens, “DISTBIC: a speaker-based segmentation for audio data indexing,” *Speech Communication*, vol. 32, pp. 111–126, 2000.
- [21] M. Cettolo, “Segmentation, classification and clustering of an Italian broadcast news corpus,” in *Conference on Content-Based Multimedia Information Access (RIAO 2000)*, Paris, France, April 2000.
- [22] M. Cettolo, M. Vescovi, and R. Rizzi, “Evaluation of BIC-based algorithms for audio segmentation,” *Computer Speech and Language*, vol. 19, pp. 147–170, 2005.
- [23] T. Kemp, M. Schmidt, M. Westphal, and A. Waibel, “Strategies for automatic segmentation of audio data,” in *Proceedings of International Conference on Acoustics Speech and Signal Processing (IEEE, ICASSP 2000)*, Istanbul, Turkey, November 2000, pp. 1423–1426.
- [24] J. Schroeder and J. Campbell, Eds., *Digital Signal Processing (DSP), a review journal - Special issue on NIST 1999 speaker recognition workshop*. Academic Press, 2000.
- [25] C. Barras and J.-L. Gauvain, “Feature and score normalization for speaker verification of cellular data,” in *Proceedings of International Conference on Acoustics Speech and Signal Processing (IEEE, ICASSP 2003)*, Hong Kong, 2003.
- [26] J. Pelecanos and S. Sridharan, “Feature warping for robust speaker verification,” in *2001: a Speaker Odyssey. The Speaker Recognition Workshop (ISCA, Odyssey 2001)*, Chania, Crete, June 2001, pp. 213–218.
- [27] D. A. Reynolds, T. F. Quatieri, and R. B. Dunn, “Speaker verification using adapted Gaussian mixture models,” *Digital Signal Processing (DSP), a review journal - Special issue on NIST 1999 speaker recognition workshop*, vol. 10, no. 1-3, pp. 19–41, 2000.
- [28] D. A. Reynolds, E. Singer, B. A. Carlson, G. C. O’Leary, J. J. McLaughlin, and M. A. Zissman, “Blind clustering of speech utterances based on speaker and language characteristics,” in *Proceedings of International Conference on Spoken Language Processing (ICSLP 98)*, 1998.
- [29] L. Nguyen, S. Abdou, M. Afify, J. Makhoul, S. Matsoukas, R. . Schwartz, B. Xiang, L. Lamel, J.-L. Gauvain, G. Adda, H. Schwenk, and F. Lefevre, “The 2004 BBN/LIMSI 10xRT English broadcast news transcription system,” in *DARPA RT04’S*, Palisades, NY, Nov 2004.
- [30] X. Zhu, C. Barras, S. Meignier, and J.-L. Gauvain, “Combining Speaker Identification and BIC for Speaker Diarization,” in *Proceedings of the 9th European Conference on Speech*

- Communication and Technology (ISCA Interspeech'05)*, Lisbon, Portugal, September 2005, pp. 2441–2444.
- [31] D. Reynolds and P. Torres-Carrasquillo, “Approaches and applications of audio diarization,” in *Proceedings of International Conference on Acoustics Speech and Signal Processing (IEEE, ICASSP 2005)*, Philadelphia, Mar 2005.
- [32] R. Sinha, S. Tranter, M. Gales, and P. Woodland, “The Cambridge University March 2005 speaker diarisation system,” in *Proceedings of the 9th European Conference on Speech Communication and Technology (ISCA Interspeech'05)*, Lisbon, Portugal, September 2005, pp. 2437–2440.
- [33] L. Canseco-Rodriguez, L. Lamel, and J.-L. Gauvain, “Speaker diarization from speech transcripts,” in *Proceedings of International Conference on Spoken Language Processing (ISCA, ICSLP 2004)*, Jeju, Oct 2004.