# Automatic Segmentation of Broadcast News Audio using Self Similarity Matrix

Sapna Soni
Institute of Technology,
Nirma University
Ahmedabad, Gujarat- 380054
Email: 12mict41@nirmauni.ac.in

Imran Ahmed
TCS Innovation labs - Mumbai,
Tata Consultancy Services Limited,
Yantra Park, Thane (West) - 400 601.
Email: ahmed.imran@tcs.com

Sunil Kumar Kopparapu
TCS Innovation labs - Mumbai,
Tata Consultancy Services Limited,
Yantra Park, Thane (West) - 400 601.
Email: sunilkumar.kopparapu@tcs.com

*Abstract*—Generally audio news broadcast on radio is composed of music, commercials, news from correspondents and recorded statements in addition to the actual news read by the newsreader. When news transcripts are available, automatic segmentation of audio news broadcast to time align the audio with the text transcription to build frugal speech corpora is essential. We address the problem of identifying segmentation in the audio news broadcast corresponding to the news read by the newsreader so that they can be mapped to the text transcripts. The existing techniques produce sub-optimal solutions when used to extract newsreader read segments. In this paper, we propose a new technique which is able to identify the acoustic change points reliably using an acoustic Self Similarity Matrix (SSM). We describe the two pass technique in detail and verify its performance on real audio news broadcast of All India Radio for different languages.

## I. INTRODUCTION

Audio segmentation is the process of partitioning an audio stream into homogeneous segments. In other words it is the process of identifying the time instants in the audio stream when there is a change in the source/speaker. Audio and speaker segmentation are necessary pre-processing step for several important tasks such as automatic transcription, indexing, summarization, speaker diarization. Segmentation is applicable to audio documents like broadcast news, telephone conversations, movies, etc. Speaker segmentation is the specific case of audio segmentation in which segments corresponding to the same speaker are identified. Speaker segmentation process generally involves the task of identifying and discarding non-speech regions in the audio; for instance silence, music, room noise, background noise or cross-talk. Several algorithms have been proposed in literature for automatic segmentation and diarization of audio data in general; and speaker segmentation in particular. An overview of different techniques is discussed in [1] and [2]. It should be noted that the techniques proposed in literature for audio segmentation are often tuned based in the type of data that they need to work on. For example, telephone conversations are spontaneous and contain frequent changes from one speaker to another while the broadcast news audio data the change in the speaker in relatively infrequent. As a result, techniques

for segmentation of broadcast news audio may not work for segmentation of telephone conversations [3] as well.

In this paper we address the problem of segmentation of broadcast news audio, with the explicit task of extracting audio segments corresponding to the newsreader (or news anchor) speech. Generally, broadcast news audio is composed of music, commercials (advertisements), news from correspondents (reporters) and recorded statements in addition to the newsreader spoken speech. Automatic extraction of the broadcast news segment corresponding to the newsreader spoken speech becomes essential for automatic time alignment of speech and the available (newsreader) text transcription. Automatic segmentation of broadcast news if of particular use in development of a frugal speech corpora using broadcast news audio (see for example [4], [5]). Techniques proposed in literature tend to face difficulties in the form of sub-optimal solutions, when used on the task of automatic segmentation and extraction of the newsreader segments. In this paper, we propose a novel technique which is able to identify the acoustic change points in an audio stream using an acoustic Self Similarity Matrix (SSM). SSM has been applied in music summarization and retrieval [6], [7] and also for audio segmentation [8], however these for a short duration audio. In broadcast audio news segmentation, the duration of the broadcast news audio is often longer than 10 mins and this results in a huge SSM, typically of the order $[60000 \times 60000]$ (assuming that we have a frame of length 10 ms) making it infeasible to extract homogeneous segments. In order to avoid this kind of computation we propose the use of a different similarity measure. We propose a two pass technique in which the long audio stream is first divided into non-overlapping segments of 5 seconds duration; and the MFCC feature vectors, extracted for every 10 ms frame is modelled as a multi-dimensional Gaussian distribution. A SSM is then computed by calculating a pairwise similarity measurement between each of these segments using Bayesian Information Criterion (BIC) [9]. The first level, coarse acoustic change points are obtained from this SSM; and then in the second pass the exact change point is obtained by computing a sliding window similarity measurement between overlapping segments of MFCC feature vectors around the identified coarse change point. We present the performance of this algorithm on

Fig. 1. Sliding window comparison technique for Speaker Segmentation.



(a) Change detection on audio news (5 change points)



(b) Change detection on audio news (6 change points)

Fig. 2. Acoustic change detection using BIC;the colored rectangles show the ground truth and the BIC based speaker segmentation obtained using LIUM.
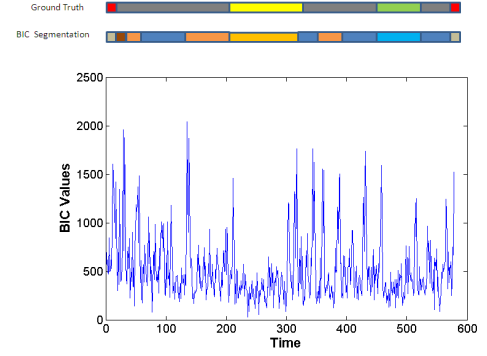
broadcast news, in different Indian languages, available from All India Radio (AIR). AIR [10] provides access to archives of news audio in several Indian languages, along with the transcripts corresponding to the newsreader spoken speech. The rest of the paper is organized as follows: Section II we briefly touch upon the issues with techniques proposed in the literature; we describe the use of SSM in Section III. We describe our approach in Section IV. Experimental results are described in Section V and we conclude in Section VI.

## II. PROBLEM WITH TRADITIONAL SPEAKER SEGMENTATION TECHNIQUES

Most of the proposed systems as described in [1], [2] consist of acoustic change detection (ACD), also called speaker change detection. ACD processes acoustic properties of an audio stream to identify instants of speaker change. The most common approach for change detection is to first divide the long audio steam into large number of smaller overlapping segments. Each small segment is represented by a set of speech feature vectors (MFCC or LSP) and then a similarity metric (BIC or KL divergence) is calculated between any two adjacent speech segments (see Figure 1). Higher the dissimilarity, the more probability that the adjacent segments belong to different speakers and hence a point of acoustic change. Typically, a chosen threshold value decides whether the segments originate from the same or a different source. Figure 2 shows the graph of the BIC (similarity) values computed every 0.1 second, between adjacent segments of 2 sec duration for two different broadcast news audio from AIR. It also includes the ground truth for segmentation, and the segmentation obtained using the LIUM - an open-source state-of-the-art toolbox for Broadcast News Diarization [11]. It can be observed that while BIC based technique works very well for one of the news audio (Figure 2 (a)), it produces over-segmentation as seen for news audio in Figure 2 (b). While popularly used, BIC based technique in literature tends to face (thresholding) problems when used on the task of automatic segmentation of the newsreader segments producing inconsistent results. We propose a new technique which identifies the acoustic change points using an acoustic Self Similarity Matrix (SSM) next.

## III. SELF SIMILARITY MATRIX (SSM) IN LITERATURE

SSM is a 2D characterization of all pairwise similarity measurements. It has been also referred to as *recurrence plots* or *dotplots*. SSM has been used in several applications including analysis of protein sequences, visualizing structure of large text corpora, detecting periodic motion in video, music segmentation and summarization (see [12] and [13]). SSM has been applied for music summarization and retrieval [6], [7] and also for audio segmentation [8]. Here music files were represented as a sequence of feature vectors and a pairwise similarity measurement between each of these feature vectors is computed to construct a SSM. In [6] the feature vectors were calculated from the audio signal at a frame rate of 20 frames per second (fps) and the typical length of the audio was around 200 seconds; whereas in [7] the feature vectors were calculated from the audio signal at a frame rate of 125 fps and the typical length of the audio was around 10 seconds. This results in the size of the SSM of $[400 \times 400]$ and $[1250 \times 1250]$ respectively. However, in the task of broadcast news audio segmentation, the broadcast news is of length more than 10 mins which results in a SSM of size $[60000 \times 60000]$ (since our task is to identify newsreader read speech, we consider a 10 msec audio segment to compute the feature vector, which results in 100 fps). In order to avoid this kind of computation we propose use a different type of similarity measurement.
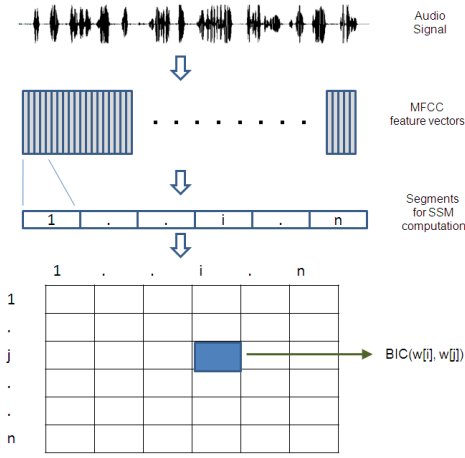
Fig. 3. Proposed technique for broadcast news segmentation.

## IV. PROPOSED TECHNIQUE FOR AUDIO SEGMENTATION

We propose a two pass technique for segmentation of broadcast news audio, in order to extract the audio segments corresponding to the newsreader (or news anchor). The first pass we identify the coarse change points in the audio stream using SSM calculated over segments of duration $2-5$ seconds and in the second pass we identify the exact change point using SSM calculated for each MFCC frame (10 msec) in the region around the (coarse) change points detected in the first pass.
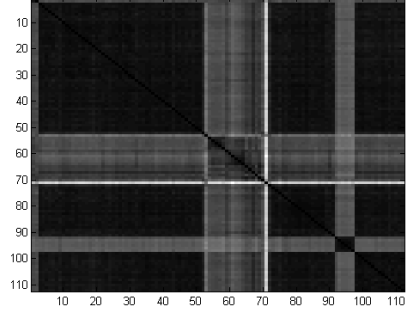
### A. First pass SSM

Figure 3 depicts the steps used for computing the SSM in the first pass. As seen, the audio signal is converted into a sequence of MFCC feature vectors (on 10 msec frames). These feature vectors are then combined together into groups of $2-5$ seconds. The grouped MFCC feature vectors in each of these segments are modelled as a multi-dimensional Gaussian distribution. A SSM is then computed by calculating a pairwise similarity measurement between each of these segments using Bayesian Information Criterion [9]. BIC is extensively used in speaker segmentation and clustering metric due to its simplicity and efficiency. The BIC similarity measure between the $i^{th}$ window ($w_i$) and the $j^{th}$ window ($w_j$) is given by
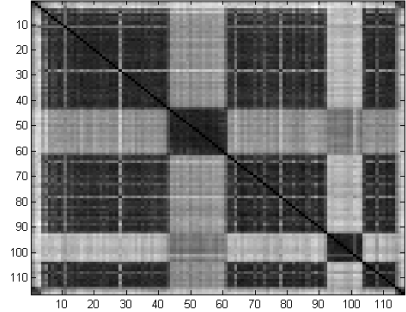
$$BIC(i,j) = \frac{N_W}{2} \log |\Sigma_W| - \frac{N_{w_i}}{2} \log |\Sigma_{w_i}| - \frac{N_{w_j}}{2} \log |\Sigma_{w_j}|,$$

where, $\Sigma_{w_i}$, $\Sigma_{w_j}$ are respectively the co-variance matrices of the feature vectors in $w_i$ and $w_j$, and $\Sigma_W$ is the co-variance matrix of all the feature vectors combined in the two windows $w_i$ and $w_j$ and $N_{w_i}/2$, $N_{w_j}/2$ and $N_W/2$ are respectively the number of feature vectors in the windows $w_i$, $w_j$ and $W$. The BIC measure is an estimate of the measure of similar between two segments; Larger values of BIC, are an indication of dissimilarity between the two segments. Figure 4 shows a visualisation of the SSM for the two news audio shown in Figure 2. The darker the regions in these SSM's the higher the similarity of the corresponding audio segments. It can

be seen from Figure 4 that as we move along the SSM in the horizontal direction the (image) edges correspond to the change in darkness represent the acoustic change points in the audio. As shown in Figure 4, it is evident that the proposed technique can reliably detect the change points in the audio stream. However, it should be noted that localization precision of the detected change point using the SSM in the first pass directly proportional to the length of the segment. Hence, the change points obtained in this pass ($2-5$ seconds) are coarse. We perform a second pass to localize and detect the actual change point.



(a) First pass SSM visualisation for audio in Figure 2(a)



(b) First pass SSM visualisation for audio in Figure 2(b)

Fig. 4. First pass SSM visualisation for audio in Figure 2

### B. Second Pass

In this pass a segment ($10-20$ second) is taken around the change point detected in the first pass and the exact change point is obtained by computing a sliding window similarity measurement between overlapping segments of MFCC features in a manner to the technique discussed in Section II. In this case it is already known that there is one change point in the entire segment and hence the highest peak in the change detection graph can be reliably selected as the change point; without having to worry about an appropriate choice of a threshold value. Once the exact acoustic change points in the audio are identified the longest segment is treated as the newsreader and other segments having a low BIC similarity with this segment are treated as the segments of the newsreader.

## V. EXPERIMENTS FOR PERFORMANCE ANALYSIS

In order to evaluate the performance of the proposed algorithm we used 10 broadcast news sourced from [10]. Each of these audio news is of 10 mins duration and consists of audio segments from the newsreader, news correspondent or reporter, commercials and music. Each of these audio files is processed to extract MFCC feature vectors, from 25 msec frames of speech, every 10 msec. In order to perform the first pass the MFCC feature vectors for each audio are grouped into a segment of 5 second duration and a SSM is computed as discussed in Section IV-A. This gives the coarse change points in the audio. Then as discussed in Section IV-B a 20 second segment is taken around the coarse change point of first pass. A sliding window BIC comparison is carried for 2 seconds of two adjacent windows, every 100 msec. The peak of the resulting BIC change detection graph is chosen as the exact change point. For the purpose of comparison the number of segments extracted using the proposed two pass algorithm is compared with the actual number of segments in the audio and also to the number of segments obtained using the LIUM open-source toolbox for Broadcast News Diarization [11].

Table I shows the results for each of the 10 news audio. It can be seen that the proposed algorithm detects almost the exact number of segments whereas the LIUM system tends to over-segment the audio files. It should be noted that two of the 10 audio news consist of only the newsreader and the proposed algorithm does not show any change points in these audio files. This can also be seen through the SSM of one these audio files, as shown in Figure 5. It should also be noted that in both the cases the actual change points were accurately marked. However, the LIUM system faces the problem of selecting change points as discussed in Section II. Whereas the proposed algorithm is able to avoid this using the SSM in the first pass.

TABLE I
COMPARISON OF NUMBER OF SEGMENTS GENERATED

| News Broadcast | Actual | by LIUM | by SSM |
|---|---|---|---|
| Hindi_Patna1Oct | 8 | 15 | 7 |
| Hindi_Indore11Oct | 7 | 21 | 7 |
| Hindi_Patna10Oct | 5 | 15 | 4 |
| Hindi_Shimla3Feb | 7 | 6 | 7 |
| Hindi_Shimla20Feb | 3 | 12 | 3 |
| Telugu_Gangtok3Oct | 1 | 1 | 1 |
| Telugu_Vijaywada3Oct | 5 | 7 | 5 |
| Telugu_Vijaywada28Sept | 5 | 4 | 5 |
| Telugu_TeluguNSD27Sept | 1 | 7 | 1 |
| Telugu_Hyderabad27Sept | 3 | 3 | 3 |
| **Total** | **45** | **91** | **43** |

## VI. CONCLUSION

In this paper we address the problem of segmentation of long audio stream like broadcast news audio, in order to extract the audio segments corresponding to the newsreader.
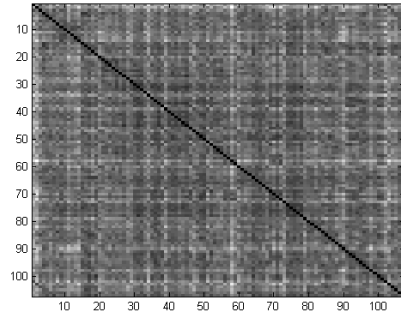


Fig. 5. SSM visualisation for broadcast news audio without any change points.

The existing systems and techniques proposed in literature give sub-optimal solutions when used for the task of automatic segmentation and extraction of the newsreader segments from a long audio stream. We proposed a two pass technique; in the first pass a number of coarse acoustic change points using an acoustic Self Similarity Matrix are identified; when in the second pass the exact position of the acoustic change points is found. The proposed algorithm is evaluated for its performance on broadcast news from All India Radio and shows to perform with higher accuracy compared to other available audio segmentation tools.

## REFERENCES

[1] S. Tranter and D. Reynolds, "An overview of automatic speaker diarization systems," *Audio, Speech, and Language Processing, IEEE Transactions on*, vol. 14, no. 5, pp. 1557–1565, 2006.

[2] X. A. Miro, "Robust speaker diarization for meetings," Ph.D. dissertation, Universitat Politecnica de Catalunya, 2006.

[3] I. Ahmed and S. K. Kopparapu, "Speaker change detection in telephone conversations," in *In Proceedings of ICSSA*, 2009.

[4] ——, "Speech recognition for resource deficient languages using frugal speech corpus," in *ICSPCC2012*, Hong Kong, China, Aug 2012, (to appear).

[5] ——, "Technique for automatic sentence level alignment of long speech and transcripts," in *In Proceedings of INTERSPEECH*, 2013.

[6] M. Cooper and J. Foote, "Summarizing popular music via structural similarity analysis," in *In Proceedings of IEEE Workshop on Applications of Signal Processing to Audio and Acoustics*, 2003.

[7] M. C. Jonathan Foote and U. Nam, "Audio retrieval by rhythmic similarity," in *In Proceedings of International Symposium on Music Information Retrieval*, 2002.

[8] J. Foote, "Automatic audio segmentation using a measure of audio novelty," in *In Proceedings of IEEE International Conference on Multimedia and Expo*, 2002.

[9] S. Chen and P. Gopalakrishnan, "Speaker, environment and channel change detection and clustering via the bayesian information criterion," in *In Proceedings of DARPA Broadcast News Transcription and Understanding Workshop*, 1998.

[10] "All india radio news archives." [Online]. Available: http://www.newsonair.com/

[11] M. R. et. al., "An opensource stateoftheart toolbox for broadcast news diarization," in *In Proceedings of INTERSPEECH*, 2013.

[12] "Self-similarity analysis: a selected bibilography." [Online]. Available: http://www.fxpal.com/?p=similaritybib

[13] "A comprehensive bibliography about rps, rqa and their applications." [Online]. Available: http://www.recurrence-plot.tk/bibliography.php