

# AN MFCC-GMM APPROACH FOR EVENT DETECTION AND CLASSIFICATION

L. Vuegen<sup>a,c,d</sup> B. Van Den Broeck<sup>b,c,d</sup> P. Karsmakers<sup>b,c,d</sup> J. F. Gemmeke<sup>a</sup> B. Vanrumste<sup>b,c,d</sup> H. Van hamme<sup>a</sup>

<sup>a</sup>ESAT-PSI, KU Leuven, Kasteelpark Arenberg 10, 3001, Heverlee, Belgium

<sup>b</sup>ESAT-SISTA, KU Leuven, Kasteelpark Arenberg 10, 3001, Heverlee, Belgium

<sup>c</sup>iMinds, Future Health Department, Kasteelpark Arenberg 10, 3001, Heverlee, Belgium

<sup>d</sup>MOBILAB, TM Kempen, Kleinhoefstraat 4, 2440, Geel, Belgium

## ABSTRACT

This abstract explores Gaussian Mixture Models (GMM) estimated from Mel Frequency Cepstral Coefficients (MFCCs) for acoustic event detection and classification. To limit the impact of silence, a shared background model is used. An average F-score of 48% for the office life subtask is obtained. However, the analysis reveals that the proposed method has difficulties to cope with the large intra-class variations (e.g. time durations, dynamic range, characteristic sounds) in the provided dataset.

**Index Terms**— Acoustic Event Detection, Mel-Frequency Cepstral Coefficients, Gaussian Mixture Models.

## 1. INTRODUCTION

The use of Gaussian Mixture Models is a well-known approach in the domain of speech- and speaker recognition applications. Research shows that this technique can achieve promising results especially in the conjunction with auditory motivated features (e.g. MFCCs) [1]. Therefore, this work will examine the use of an MFCC-GMM baseline acoustic event detector and classifier on the publicly available database from the IEEE-AASP challenge. The remainder of this abstract is organized as follows: feature extraction and training phase will be briefly discussed in section 2. Section 3 handles about the used event detector and classifier. The executed experiments and obtained results from the subtasks office life and office synthetic are given in section 4. Finally, the conclusions are discussed in section 5.

## 2. FEATURE EXTRACTION AND TRAINING

Figure 1 is a flowchart of the feature extraction and training phase which has been used during this work. This process starts by iteratively loading the waveform (.wav) files from each event. Next, the corresponding MFCC features, including the first and second derivatives, are computed.

Labeling the extracted features into the actual event features and background features happens in two stages. First, the provided annotation files from the 2 different annotators are used to locate the event features. The earliest onset mark of both annotators is used as onset and the latest offset mark of both is used as offset in order to reduce the probability of labeling event features as belonging to the background. Next, a threshold on the first MFCC-coefficient, further denoted as  $C_0$ , is applied to remove the within event silences (e.g. silence between 2 phone rings) from the remaining event features. Frames with a  $C_0$  lower than a threshold can be assumed as low-energetic and are therefore added to the background features.

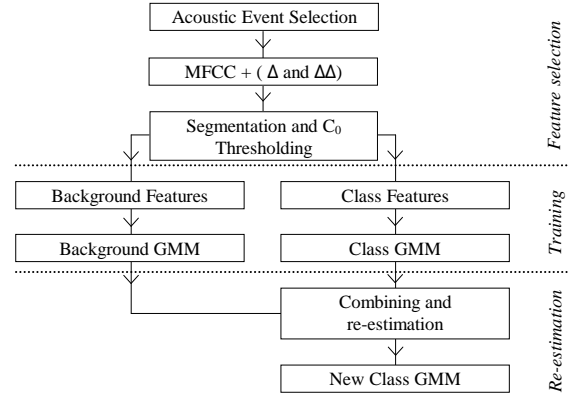


Figure 1: Flowchart of feature extraction and training phase.

The training process starts by estimating a shared background GMM and all the class GMMs (16 in total) on basis of the background features and class event features respectively. This by applying the Expectation-Maximization (EM) algorithm as explained in [1].

Finally, the class GMMs will be combined and re-estimated in the presence of the background model which is not re-estimated. This is preferred over relying on each of the class-GMMs to model the silence frames independently. This way, the shared background GMM will produce the same score for each of the class assumptions and hence the impact of silence frames on the model likelihoods will be minimized. Reestimation of the  $i$ -th Gaussian in a class mixture is achieved by replacing its posterior with:

$$p(i|x, \lambda) = \frac{\lambda_i N(x_i | \mu_i, \Sigma_i)}{\sum_{j=1}^M \lambda_j N(x_i | \mu_j, \Sigma_j) + \lambda_{prior} \sum_{k=1}^N \lambda_{back,j} N(x_i | \mu_{back,j}, \Sigma_{back,j})}, \quad (1)$$

Compared to the standard EM-algorithm is (1) expanded with an additional term in the denominator, i.e. the contribution of the background GMM to the data likelihood. The weighting prior defines the amount of probability mass that is assigned during the maximization step to the background model and is determined in (2). During the initialization of the EM-algorithm it is required to set an initial value for  $\lambda_{prior}$  because the sum over the class weights still unities (property of GMMs).

$$\lambda_{prior} = 1 - \left( \sum_{j=1}^M \lambda_j \right) \quad (2)$$

In the experimental setup will the influence of the proposed technique examined by combining the shared background GMM with the class GMMs. This by a) applying the adapted EM-formula with an initial  $\lambda_{prior}$  of 0.2 and b) a linear combination of the background and class GMMs with a ratio 1/5 respectively.

### 3. EVENT DETECTOR AND CLASSIFIER

The event detector and classifier for the subtask office life starts by extracting the MFCC features (including 1<sup>st</sup> and 2<sup>nd</sup> derivative) from the acoustic event script (see Figure 2). A posteriorgram is computed by comparing these features with both the estimated background model and all the class GMMs. Next, the posteriors from each class are moving averaged filtered with a window size depending on the minimum class duration observed in the training dataset. This smoothens the class posteriors and takes the minimum occurring time duration of each class more or less into account.

Detecting events in the office life subtask is based on  $C_0$  thresholding. It can be assumed that an event has occurred when the value of  $C_0$  was above a predefined threshold during a certain period of time. The values of  $C_0$  and minimum time duration are defined experimentally and further mentioned in the next section.

As last comes classifying the detected events. This is simply done by determining which GMM model produces the highest averaged a-posteriori score. In case that the detected event is classified as background it will be neglected and therefore removing it as an occurred event.

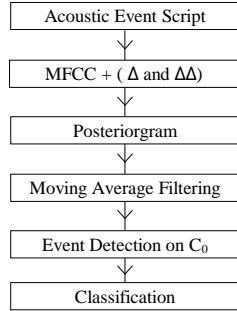


Figure 2: Flowchart of the event detector and classifier.

### 4. EXPERIMENTS AND RESULTS

In order to determine the performance of the proposed method are following parameters examined in function of the averaged event based F-score: a) resampling to a lower sample frequency i.e. 16kHz, b) influence of the first and second derivative, c) the number of Gaussians and d) re-estimation or just linear combination both with a  $\lambda_{prior}$  of 0.2.

During this experiment, the minimum value of  $C_0$  was set to -150 and -189.5 (determined experimentally) for a sample frequency of 44,1kHz and 16kHz respectively. Also the window sizes of the moving average filter applied on the class posteriors was set as half of the minimum event duration of each class. Figure 3 shows the obtained averaged F-scores and following observations are made:

- Applying a down sampling to 16kHz has for the most parameter combinations a positive effect on the F-score. A possible explanation is that the higher frequency bands contain more noise than actual characteristic information of the occurring event.
- The usage of the first and second derivative has a small positive effect on the F-score.
- Applying the proposed re-estimation algorithm does not increase the F-score.

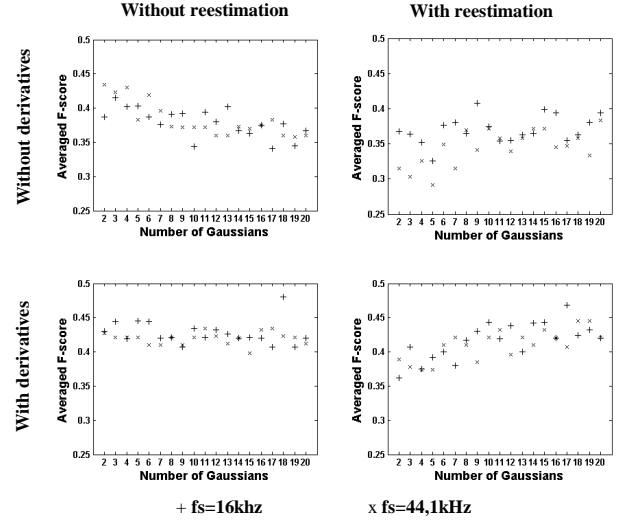


Figure 3: Averaged results of the event based F-score

Table 1 on the other hand gives the associated evaluation metrics corresponding to the highest achieved averaged F-score (in Figure 3). i.e. 46,9% and 48% for with and without re-estimation respectively. These F-scores occurs both when the derivatives are included and a resampling to 16kHz is applied. The corresponding number of Gaussians are 17 (with re-estimation) and 18 (without re-estimation) and as one can see, no major differences occurs between both methods.

Table 1: Results on the Office Live Dataset for various metrics.

Metric	Evaluation Method					
	Event Based		Class-Wise Event Based		Frame Based	
Re-estimation	Y	N	Y	N	Y	N
R	36,4	37,4	39,9	40,9	33,7	38,2
P	68,2	69,6	44,0	40,1	81,7	84,6
<b>F-score</b>	<b>46,9</b>	<b>48</b>	<b>37,8</b>	<b>38,2</b>	<b>50,3</b>	<b>52,2</b>
AEER	0,99	0,96	0,93	0,86	0,80	0,76
Offset R	30,0	30,1	31,1	32,2	-	-
Offset P	56,3	58,3	37,7	35,0	-	-
Offset F-score	38,6	39,9	31,0	31,2	-	-
Offset AEER	1,89	1,17	1,19	1,12	-	-

Table 2 and 3 gives the corresponding F-scores of the office synthetic task. The same parameters were used as in Table 1 however the minimum  $C_0$  was changed to a value just above the noise floor, i.e. -35, -45 and -95 for the SNR of -6, 0 and 6 respectively.

As one can see, the achieved results are dropped extremely, even for the easiest combination, i.e. a SNR of 6 and the lowest degree of overlapping. One of the reasons of a lower score is because our detection algorithm expects only 1 event when an event is detected during a certain time span. Second, research shows that GMMs are extremely independent to the contribution of noise. Even the smallest amount of noise can cause an enormous drop of performance [4].

Table 2: F-scores on the Office Synthetic Dataset (with re-estimation).

		Density	SNR		
			-6	0	6
Evaluation method	Event based	low	16,7	0	0
		medium	0	0	9,1
		high	0	1,87	2,08
	Class-wise event Based	low	16,7	0	0
		medium	0	0	5,56
		high	0	2,22	2,22
	Frame Based	low	12,6	0	0
		medium	0	0	5,18
		high	0,43	2,59	5,24

Table 3: F-scores on the Office Synthetic Dataset (without re-estimation).

		Density	SNR		
			-6	0	6
Evaluation method	Event based	low	16,7	0	0
		medium	0	0	9,09
		high	0	1,87	2,08
	Class-wise event Based	low	16,7	0	0
		medium	0	0	5,4
		high	0	2,22	2,67
	Frame Based	low	12,6	0	0
		medium	0	0	5,18
		high	0,43	1,09	5,98

## 5. DISCUSSION

The overall performance of the proposed method were not so promising as hoped, especially for office synthetic subtask. The most obvious explanation is that Gaussian Mixture Models have difficulties to cope with a) the large variation in the characteristic sounds of some classes (e.g. phone and alert) and b) the relative low amount of training examples in the dataset. This results in a harder classification problem and therefore reducing the accuracy of the classifier. Besides, the large variation in time duration and energy increases the difficulty of the detection task and therefore also decreasing overall the performance of the system.

## 6. ACKNOWLEDGMENT

This work was performed in the context of following projects: ALADIN (IWT-SBO project contract 100049), IWT doctoral scholarships (contract 111433 and 121565) and FallRisk. The iMinds FallRisk project is cofounded by iMinds (Interdisciplinary Institute for Technology), a research institute founded by the Flemish Government. Companies and organizations involved in the project are COMmeto, Televic Healthcare, TP Vision, Verhaert and Wit-Gele Kruis Limburg, with project support of IWT.

## 7. REFERENCES

- [1] D. A. Reynolds and R.C. Rose, "Robust Text-Independent Speaker Identification Using Gaussian Mixture Speaker Models", in *IEEE Transactions on Speech and Audio Processing*, Vol. 3, No. 1, pp. 72–83, 1995.
- [2] A. Mesaros, T. Heittola, A. Eronen, And T. Virtanen, "Acoustic Event Detection In Real Life Recordings", in *the 2010 European Signal Processing Conference*, 2010.
- [3] S. Ntalampiras, "A Novel Holistic Modelling Approach for Generalized Sound Recognition", in *IEEE Signal Processing Letters*, Vol. 20, No. 2, February 2013.
- [4] L. Vuegen, P. Karsmakers And B. Vanrumste, "Comparative Study of Sound Recognition Algorithms", in *K.H. Kempen proceedings of Master Theses*, June 2012.