



# Automatic Summarization of Cricket Highlights using Audio Processing

Ritwik Baranwal

Information Technology, Maharaja Agrasen Institute of Technology, New Delhi, India

## To Cite this Article

Ritwik Baranwal, "Automatic Summarization of Cricket Highlights using Audio Processing", *International Journal for Modern Trends in Science and Technology*, Vol. 07, Issue 01, January 2021, pp.- 48-53.

## Article Info

Received on 22-November-2020, Revised on 18-December-2020, Accepted on 22-December-2020, Published on 29-December-2020.

## ABSTRACT

*The problem of automatic excitement detection in cricket videos is considered and applied for highlight generation. This paper focuses on detecting exciting events in video using complementary information from the audio and video domains. First, a method of audio and video elements separation is proposed. Thereafter, the "level-of-excitement" is measured using features such as amplitude, and spectral center of gravity extracted from the commentators speech's amplitude to decide the threshold. Our experiments using actual cricket videos show that these features are well correlated with human assessment of excitability. Finally, audio/video information is fused according to time-order scenes which has "excitability" in order to generate highlights of cricket. The techniques described in this paper are generic and applicable to a variety of topic and video/acoustic domains.*

**KEYWORDS:** Video Segmentation, Audio Chunks, Short Time Energy.

## I. INTRODUCTION

This study focuses on the problem of identifying exciting-events in multimedia content. Our approach analyzes speech characteristics that identify islands (or "hot-spots") of strong emotion. In general, the ability to automatically parse multimedia content and tag "interesting events" is important for many domains such as sports, security, movies/TV shows, broadcast news, etc. A number of technologies such as search, summation, and mash-ups, can utilize "hot-spot" information to enhance access to, as well as navigation of content. For example, emotional "hot-spots" within sports videos are very likely to be "exciting" and this information can be used to guide the process of automatically generating highlights. This constitutes the motivation for this work, where automatic highlights of cricket videos

are generated using emotional "hot-spot" detection (or "exciting events" detection).

Researchers have utilized audio and video streams to extract features that identify exciting plays in sports videos. Among video-based features, motion and density of cuts have been found to be useful for detection[1]. On the other hand, audio-based features have been derived from both speech (generally commentators) and background (generally audience), where audience-events like cheering/applause as well as the commentators speech characteristics have proven to be useful [2,3]. While video-based features tend to be more game-dependent, audio-based feature detecting exciting plays. Research in audio-based features have focused on emotion analysis of the commentator's speech and

employ this information with heuristics to identify exciting plays.

In this paper, we present a novel approach for auto-curating sports highlights, showcasing its application for cricket match. Our approach combines information from the player, spectators, and the commentator to determine a game's most exciting moments. A common trait among most of the sports is that whenever an exciting moment is happening the commentators speak loudly and the crowd cheers. We have used this as a signature for finding those moments where important or exciting things are happening in the match. So we just needed to analyze the audio of the match find those moments where the crowd cheered or commentators are excited and extracted those parts.

## **METHOD**

In our approach the generation is done on basic analysis of sound or audio processing, we know that whenever an important event occurs during a match the adrenaline rush can be seen into commentary, so in simple terms it's an impulsive energy in short span of time. So, what we have done here is we have extracted all those impulsive commentary events and have seen the corresponding video at the same timestamp and generated the desired highlights.

## **STRUCTURE OF PAPER**

The paper is organized as follows: In Section 1, the introduction of the paper is provided along with the structure, important terms, objectives and overall description. In Section 2 we discuss related work. In Section 3 we have the complete information about image processing tools. Section 4 shares information about the flexible YAML templating system created for it, its advantages and disadvantages. Section 5 tells us about the methodology and the process description. Section 6 tells us about the future scope and concludes the paper with acknowledgement and references.

## **OBJECTIVES**

This project aims to address some of the problems in current systems by greatly minimizing the human intervention in the process and thus reducing costs and errors. The aim is to ease the task of both the technicians and audience.

## **II. RELATED WORK**

### **Video Summarization.**

There is a long history of research on video summarization [4], which aims at producing short videos or keyframes that summarize the main content of long full-length videos, by looking at eliminating redundancy either at signal level (feature dimensionality reduction [5]) or in semantic content [6]. Our work also aims at summarizing video content, but instead of optimizing for representativeness and diversity, as traditional video summarization methods do, our goal is to find highlights or exciting moments in the videos. A few recent methods address the problem of highlight detection in consumer videos [7]. Instead our focus is on sports videos, which offer more structure and objective metrics than unconstrained consumer videos.

### **Automatic Trailer Generation.**

Another sub-area of video summarization involving multimodal video analysis that goes beyond content recognition, and focusing instead on affective responses evoked by the video is movie trailer generation [8,9]. For example, Evangelopoulos et al. [9] model and combine audio visual and textual saliency to select the most relevant scenes in a movie. In this space, works focus on detecting content with the highest emotional impact based on movie genre. For instance, in horror movies scenes evoking feelings of suspense or fear are important [9]. In our domain of interest, on the other hand, only positive emotions connected to excitement are relevant. Further-more, differently from this line of research, the focus of our work is on identifying and measuring subjects reactions (players, crowd, and commentator)

directly in the video stream, rather than inferring reactions which are supposed to be evoked by inspected content which is deemed as "impressive" [1].

### **Sports Highlights Generation.**

Several methods have been proposed to automatically extract highlights from sports videos based on audio and visual cues. Example approaches include the analysis of replays [1, 2, 4], crowd cheering [6, 3], motion features [5], and closed captioning [4]. More recently, Bettadapura et al. used contextual cues from the environment to understand the excitement levels within a basketball game. Tang and Boring [3] proposed to automatically



produce high-lights by analyzing social media services such as twitter. Decroos et al. [8] developed a method for forecasting sports highlights to achieve more effective coverage of multiple games happening at the same time. Different from existing methods, our proposed approach offers a unique combination of excitement measures extracted from live video streams to produce highlights, including information from the spectators, the commentator, and the player reaction. As such, our system incorporates and combines most of the information employed by previous works (audio, visual, text).

It could also be easily extended to integrate other sources of attention or excitement, such as social media feeds or production cues (replays, closed captions, etc.). In addition, we enable personalized highlight generation or retrieval based on a viewer's favorite players.

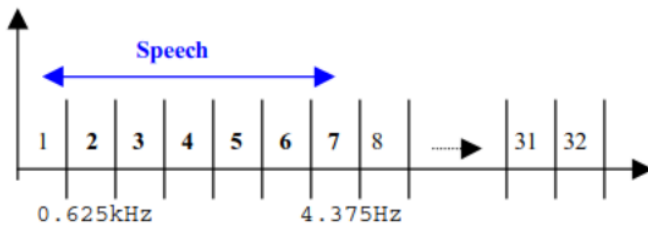
### **Self-Supervised Learning.**

In recent years, there has been significant interest in methods that learn deep neural network classifiers without requiring a large amount of manually annotated training examples. In particular, self-supervised learning approaches rely on auxiliary tasks for feature learning, leveraging sources of supervision that are usually available "for free" and in large quantities to regularize deep neural network models. Examples of auxiliary tasks include the prediction of ego-motion [4, 1], location and weather [4], spatial context or patch layout [2, 4], image colorization [4], and temporal coherency [2]. Altar. [5] explored the natural synchronization between vision and sound to learn an acoustic representation from unlabeled video. We leverage this work to build audio models for crowd cheering and commentator excitement using few training examples, and use those classifiers to constrain the training data collection for player reaction recognition. More interestingly, we exploit the detection of TV graphics as a free supervisory signal to learn feature representations for player recognition from unlabeled video.

### **III. AUDIO PROFILE GENERATION**

a) MPEG Bitstream Processing The Fischlár system captures television broadcasts and encodes the programmes according to the MPEG-1 digital video standard with the audio signal coded in line with the Layer-II profile [8]. Unlike many other audio compression algorithms, which make assumptions about the nature of the audio source, MPEG-1 Audio exploits the perceptual restrictions of the human auditory system, via psychoacoustic weighting of the bit allocation for each frequency subband, to attain its compression [7]. The MPEG-1 Layer-II compression algorithm encodes audio signals by dividing the frequency spectrum of the audio signal, bandlimited to 20 kHz into 32 subbands which approximate the ear's critical bands. The subbands are assigned individual bit-allocations according to the audibility of quantisation noise within each subband. A psychoacoustic model of the ear analyses the audio signal and provides this information to the quantiser. Layer-II frames consist of 1152 samples; 3 groups of 12 samples from each of 32 subbands. A group of 12 samples gets a bit-allocation and, if this is non-zero, a scale factor. Scale factors are weights that scale groups of 12 samples such that they fully use the range of the quantiser (the encoder uses a different scale factor for each of the three groups of 12 samples within each subband only if necessary). The scale factor for such a group is determined by the next largest value (given in a look-up table) to the maximum of the absolute values of the 12 samples. Thus it provides an indication of the maximum power exhibited by any one of the 12 samples within the group [9, 1].

b) Amplitude of the Speech Band Most of the energy in a speech signal lies between 0.1 kHz – 4 kHz. According to the MPEG-1 Layer-II audio standard, the maximum allowable frequency component in the audio signal is at 20 kHz. At the encoder, the frequency spectrum (0 – 20 kHz) is divided uniformly into 32 subbands, each having a bandwidth of 0.625 kHz [4]. Thus, subbands 2 through 7 represent the frequency range from 0.625 kHz – 4.375 kHz. See Figure 1. Speech 1 2 3 4 5 6 7 8 31 32 0.625 kHz 4.375 Hz



**Figure 1: MPEG-1 Layer-II Frequency Subbands**

Figure 1: MPEG-1 Layer-II Frequency Subbands For sports programme audio tracks, by strictly limiting the audio examination to these subbands, which approximate the range of the speech band, we further concentrate the audio investigation on commentator vocals. Therefore, the influence of the commentator on the generation of the audio amplitude profile is increased. This is clearly desirable since it bolsters the assumption that the profile will be an accurate indicator of the significance of the content. It was expected that the examination of subbands 2 through 7 would provide for a reasonable trade-off between rejection of low-frequency background noise (typically present in sports programmes which would naturally upset results) and the capture of the fundamental frequency for excited speech.

c) Boundary Detection One of the problems with the audio amplitude technique is caused by the inclusion of supplementary content which typically accompanies the main event in a sports programme. Features such as player profiles, highlights of recent events etc. tend to contain attributes such as commentator dialogue and spectator noise, similar to that of the main event. The problem is that these features generally have audio amplitudes comparable to that of the event of interest. To combat this problem, the system must be able to detect the temporal boundaries of the main feature within the overall sports programme. This is done by searching through the audio track for extended periods of sustained volume. Segments such as interviews, studio discussions, archive video clips, etc. which make up the peripheral content, are flagged by the intermittent occurrence of brief moments of silence. For example short silences exist in between sentences spoken by an anchor person, when switching from anchor person to video clips, or between advertisements. In contrast, the main event in a sports broadcast features relatively long periods of sustained volume due to the continuous presence of background noise. On this basis it may be automatically distinguished from the supplementary content. i.e.

the temporal boundaries of the main event within the overall program may be detected. For the summary generation, the probing domain is restricted to lie within these boundaries.

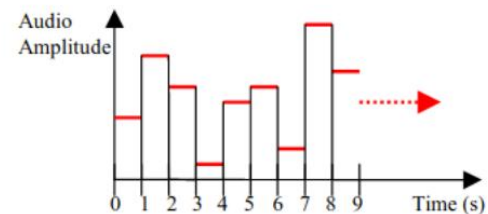
#### IV. CASE STUDY

##### a) Task

The following is an illustration of the automatic generation of a 10-minute summary of a terrestrial TV broadcast of a sports event via the discussed technique. The experimental subject is the UEFA Cup Final 2001 featuring Liverpool FC Vs Alaves FC. This was a near 3-hour soccer match broadcast, resulting in a 5-4 victory for Liverpool FC. The programme featured the main event plus studio discussions and analysis, player profiles, highlights of related events and advertisement breaks.

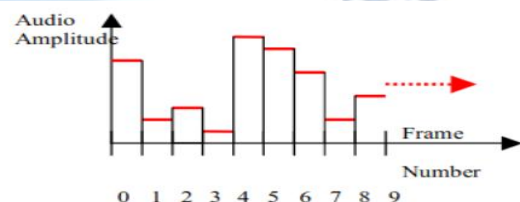
##### b) Amplitude Profiles

A second-by-second audio amplitude profile was established by a superposition of all the scale factors from subbands 2-7 over a window length of one second. See Figure 2.



**Figure 2: Per-Second Audio Amplitude Profile**

A frame-by-frame audio amplitude profile was established by a superposition of all the scale factors from subbands 2-7 over a window of length corresponding to one video frame ( $\approx 1/25$ s). See Figure 3.



**Figure 3: Per-Frame Audio Amplitude Profile**

##### c) Boundary Detection

The overall structure of the near 3-hour subject, as captured by Físchlár, is described below. In terms of summary generation, segments of interest are identified by an asterisk.



\*1st\_half~51mins  
 Studio\_analysis~14mins  
 \*2nd\_half~49mins  
 Studio analysis~4 mins  
 \*Extra\_time~26mins  
 Studio analysis~6 mins

A silence threshold was empirically determined as

**Sth = 0.033 \* overall mean audio amplitude**

Using the per-frame audio amplitude profile and Sth, the entire audio track of the subject was examined for periods of continuous volume lasting, at least 1-minute. It was found that sustained volumes exceeding Sth occur during the following videoframes:

309 – 3504~2mins  
 3867 – 6608~1min  
 6984 – 13577~4mins  
 15037 – 19467~3mins  
 19553 – 23405~2mins  
 \*26248 – 102751~51mins  
 106225 – 109935~2mins  
 112696 – 115354~1min  
 \*123629 – 198950~50mins  
 199586 – 201168~1min  
 \*201252 – 244086~28mins  
 245527 – 247690~1min

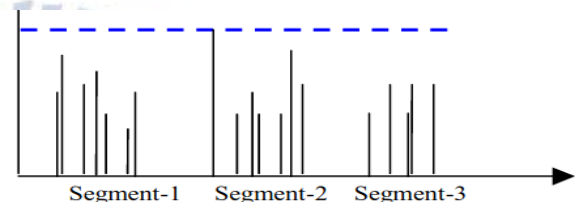
Further thresholding at a length of 10-minutes rejects all segments except for three (identified by an asterisk), which correspond almost precisely to these segments of interest mentioned previously (i.e. the temporal boundaries of the match play segments were accurately detected). Changing units to seconds these are:

- Segment-1:1050s – 4110s
- Segment-2 :4945s – 7958s
- Segment-3:8050s – 9763s

Only the content which resides within these boundaries is eligible for inclusion in the summary. Hence, further audio processing is restricted accordingly. The boundary detection preamble is not a crucial component of the summarisation procedure, i.e. in the event of failure, the main audio analysis procedure would still be expected to produce a moderately successful summary. However, it is a beneficial tool which prevents the consideration of irrelevant material

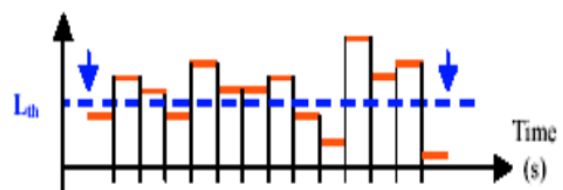
and thus lightens the workload of subsequent procedures.

d) Summary Generation The per-second audio amplitude profiles of segments 1-3 (above) were examined. A loudness threshold,  $L_{th}$ , was defined and initialised to the value corresponding to the largest peak found. See Figure 4.



**Figure 4:** Examination of Segments 1-3

An audio amplitude peak is defined as loud if it exceeds  $L_{th}$ . Ignoring isolated peaks,  $L_{th}$  was gradually reduced until it began to pick out loud periods of at least 3-seconds in duration (audio surges). See Figure 5.



**Figure 5:** Decreasing  $L_{th}$  and detecting audio surges

Figure 5 shows three sections which extend beyond the current value of  $L_{th}$ . The second and third have time spans of 4 seconds and 3 seconds respectively. Thus both are recognised as audio surges. The first section is ignored since with a length of 2 seconds, it does not meet the minimum surge threshold of 3-seconds.  $L_{th}$  was further reduced until the amount of detected surges was sufficient such that a 10-minute video summary could be produced. The summary was then generated by first matching up the video clips within the combined audio/video track which temporally align with the audio surges. Then, a pre-clip buffer of 1 shot and a post clip buffer of 2 shots was appended (to make viewing the amalgamation less visually disturbing). Finally these clips were extracted from the audio/video stream and (chronologically) concatenated to generate the highlights summary. See Figure 6.

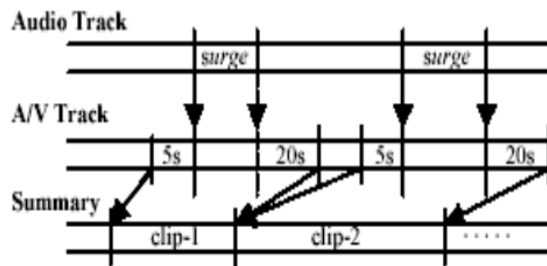


Figure 6: Summary generation

## V. RESULTS

The analysis returned 18 individual clips corresponding to the following descriptions, comprising a summary length of just over 5-minutes:

1. Wicket \*
2. Boundary -
3. Boundary -
4. Wicket -
5. Wicket -
6. Boundary \*
7. Boundary -
8. Boundary -
9. Boundary -
10. Boundary \*
11. Boundary \*
12. Boundary #
13. Wicket #
14. Boundary #
15. Catch Drop -
16. Boundary -
17. Wicket -
18. Boundary -

For the purposes of evaluation, the nineteen clips returned were examined and classified into four categories according to significance. Twelve clips seemed to depict very significant moments of the feature and hence were described as definite highlights (-). The inclusion of definite highlights in the summary is always preferred. Four of the

clips returned seemed to represent moments of arguably lesser significance. These were described by the term quasi-highlights (#), and their inclusion in the summary is desired once all definite highlights already have been. The system returned three further clips containing content of considerably less significance, labeled lowlights (\*). Inclusion of lowlights would typically not be tolerated except when the combined length of all definite and quasi-highlight clips fails to satisfy the desired length of the summary.

## VI. FUTURE SCOPE AND CONCLUSION

Further this approach can be used with OCR to increase the percentage of definite highlights. Different threshold combinations can also be tried to filter out better results. In this study, a novel methodology that uses estimates of excitability in sports video to create automatic highlights was presented. First, a method of audio and video elements separation is proposed. Thereafter, the “level-of-excitement” is measured using features such as amplitude, and spectral center of gravity extracted from the commentators speech’s amplitude.

## REFERENCES

- [1] C.Liu, Q. Huang, S. Jiang, L. Xing, Q. Ye, and W. Gao, “A frame-work for flexible summarization of racquet sports video using mul-tiple modalities,” *Computer Vision and Image Understanding*, vol. 113, pp. 415–424, 2009.
- [2] E. Kijak, G. Gravier, P. Gros, L. Oisel, and F. Bimbot, “HMM based structuring of tennis videos using visual and audio cues,” in *ICME*, 2003.
- [3] R. Radhakrishnan, Z. Xiong, A. Divakaran, and Y. Ishikawa, “Gen-eration of sports highlights using a combination of supervised and unsupervised learning in audio domain,” in *Pacific Rim Conference on Multimedia*, 2003.
- [4] Y.-F. Ma, L. Lu, H.-J. Zhang, and M. Li, “A user attention model for video summarization,” in *ACM Multimedia*, 2002.
- [5] J. Zhang, J. Yu, and D. Tao, “Local deep-feature alignment for unsupervised dimension reduction,” *IEEE Transactions on Image Processing*, vol. 27, no. 5, pp. 2420–2432, 2018.
- [6] K. Zhang, W.-L. Chao, F. Sha, and K. Grauman, “Video summarization with long short-term memory,” in *ECCV*, 2016.
- [7] M. Sun, A. Farhadi, and S. Seitz, “Ranking domain-specific highlights by analyzing edited videos,” in *ECCV*, 2014.
- [8] G. Irie, T. Satou, A. Kojima, T. Yamasaki, and K. Aizawa, “Automatic trailer generation,” in *ACM Multimedia*, 2010.
- [9] G. Evangelopoulos, A. Zlatintsi, A. Potamianos, P. Maragos, K. Rapantzikos, G. Skoumas, and Y. Avrithis, “Mul-timodal saliency and fusion for movie summarization based on aural, visual, and textual attention,” *IEEE Transactions on Multimedia*, vol. 15, no. 7, pp. 1553–1568, 2013.