

A Hierarchical Framework for Semantic Scene Classification in Soccer Sports Video

M. H. Kolekar

Department of Computer Science,
University of Missouri,
Columbia, USA
Email: kolekarm@missouri.edu

K. Palaniappan

Department of Computer Science,
University of Missouri,
Columbia, USA
Email: palaniappank@missouri.edu

Abstract— In this paper, we propose a novel hierarchical framework for soccer (football) video classification. Unlike most existing video classification approaches, which focus on shot detection followed by classification based on clustering using shot aggregation, the proposed scheme perform a top-down video scene classification which avoids shot clustering. This improves the classification accuracy and also maintains the temporal order of shots. In the hierarchy, at level-1, we use audio features, to extract potentially interesting clips from the video. At level-2, we classify these clips into field view and non-field view using feature of dominant grass color ratio. At level-3a, we classify field view into three kinds of views using motion-mask. At level-3b, we classify non-field view into close-up and crowd using skin color information. At level-4, we classify close-ups into the four frequently occurring classes such as player of team-A, player of team-B, goalkeeper of team-A, goalkeeper of team-B using jersey color information. We show promising results, with correctly classified soccer scenes, enabling structural and temporal analysis, such as highlight extraction, and video skimming.

I. INTRODUCTION

In recent years, sports video has emerged as important area of research due to its wide viewership, ease of digital archival and huge commercial potential [1], [2]. Moreover, the distribution of sports video across the Internet further increases the need for automatic video analysis, such as quick browsing, video summarization, detecting and recording interesting highlights for later review. Automatic video analysis remains a challenging area due to the large gap between semantics and low level feature interpretation of video data, which has led to domain specific approaches.

Sports video analysis has received much attention in the area of digital video processing. Existing approaches can be broadly classified as genre-specific or genre-independent. Due to dramatically distinct broadcast styles for different sports genres, much of the prior art concerns genre specific approaches. Researchers have targeted the individual sports game such as soccer (football) [3], tennis [4], cricket [5], basketball [6], volleyball [7], etc. These works show that genre specific approaches typically yield successful results within the targeted domain. In comparison with the genre-specific research work, less work is observed for genre-independent studies [8], [7]. For a specific sports event detection task, it is not feasible to expect a general solution that will work

successfully across all genres of sports video.

In the case of soccer sports video, Li et. al. [9] proposed rule based algorithm using low-level audio/video features for soccer video summarization. Babaguchi et. al. [10] proposed event detection by recognizing the textual overlays from soccer video. Wan et. al. [11] proposed dominant speech features to generate soccer highlights. Ding et. al. [12] proposed segmental Hidden Markov Model for view-based soccer analysis. Wang et. al. [3] proposed an automatic approach for personalized sports music video generation.

Most of the research in sports video processing [13], [14] assumes a temporal decomposition of video into its structural units such as scenes, shots and frames similar to other video domain including television and films. A shot refers to a group of sequential frames often based on single set of fixed or smoothly varying camera parameters (i.e. close-up, medium or long shots, dolly, pan, zoom, etc). A scene refers to a collection of related shots. In our work, we extract the scenes and after analysis assign a descriptive label (i. e. event) to each scene. The main contributions of the paper are as follows: 1) We propose novel hierarchical framework for soccer video, 2) Our proposed field-view classification based on motion-mask is new and gives very good classification accuracy, 3) We propose novel domain specific close-up detection and classification algorithm based on skin detection and jersey color comparison. The rest of the paper is organized as follows. Section-II presents proposed hierarchical classification tree. Section-III presents experimental results. Section-IV concludes the paper with direction for future work.

II. HIERARCHICAL CLASSIFICATION

In [3] the authors integrate multiple features to classify video sequences. Although the integration of multiple features improves the classification accuracy, it leads to other problems such as proper selection of features, proper fusion and synchronization of right modalities, critical choice of the weighting factor for the features and computational burden. To cope with these problems, we propose a novel hierarchical classification framework for the football videos as shown in Figure 1, which has the following advantages: (1) The

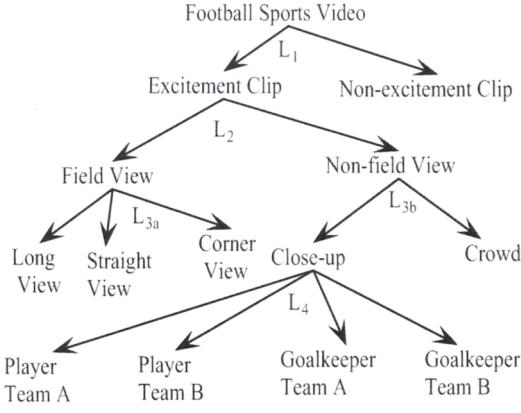


Fig. 1. Tree Diagram of Hierarchical Framework

approach avoids shot detection and clustering that are the necessary steps in most of video classification schemes, so that the classification performance is improved. (2) The approach uses top-down four-level hierarchical method so that each level can use simple features to classify the videos. (3) This improves the computation speed, since the number of frames to be processed will remarkably reduce level by level.

A. Level-1: Excitement Detection

We have observed the following changes in sports audio track when exciting events occur: 1) Spectator's cheer and commentator's speech becomes louder; 2) Commentator's talk becomes more rapid. Based on these observations, we have used two popular audio content analysis techniques- short-time audio energy and zero crossing rate (ZCR) for extracting excitement clip. We are considering the short-time as the number of audio samples corresponding to one video frames. A particular video frame is considered as an excitement frame if the product of its audio excitement and ZCR exceeds a certain threshold. After computing short-time audio energy $E(n)$ and ZCR $Z(n)$ [15], we have proposed following steps for excitement clip detection.

Averaging through sliding window:

To distinguish genuine audio excitement from audio noise, we have used a sliding window. However, it helps for early detection of the events as well.

$$E'(n) = \frac{1}{L} \sum_{l=0}^{L-1} E(n+l) \quad \text{and} \quad Z'(n) = \frac{1}{L} \sum_{l=0}^{L-1} Z(n+l) \quad (1)$$

where, L is the length of sliding window.

The normalized values are as follows:

$$E''(n) = \frac{E'(n)}{\max_{1 \leq i \leq N} E'(i)} \quad \text{and} \quad Z''(n) = \frac{Z'(n)}{\max_{1 \leq i \leq N} Z'(i)} \quad (2)$$

where, N is the total number of video frames.

Excitement frame detection

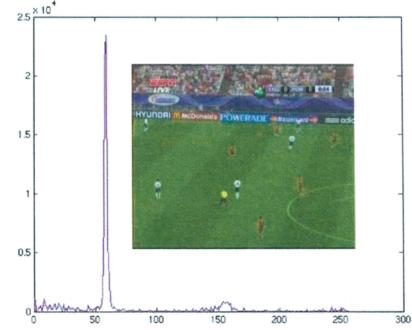


Fig. 2. Histogram of hue component of field view image, which is shown in the figure

The product $P(n)$ is given as

$$P(n) = E''(n) * Z''(n) \quad (3)$$

Based on the product term $P(n)$, a video frame n will be finally labeled as $\psi(n) \in [0, 1]$ as defined below:

$$\psi(n) = \begin{cases} 1 & (\text{excitement}) \quad P(n) \geq \mu_p \\ 0 & (\text{non-excitement}) \quad \text{otherwise} \end{cases} \quad (4)$$

where, μ_p is the mean of $P(n)$.

Excitement clip detection

To distinguish genuine audio excitement from audio noise, we select the excitement clips of duration greater than 20 seconds.

B. Level-2: Field View Detection

At level-2, we are using grass-pixel ratio similar to [16] to classify the excitement clips into field view and non-field view. In our experimental set-up, we consider 70 field view images in hsv format for training. We plot 256-bin histogram of the hue component of these images. We pick up the peaks of hue histogram of these images. As shown in Figure 2, we observed peak at bin $k = 59$ and value of the peak is 23486 for the particular image of size 288 by 352. By testing all 70 images, we observed that the green color peak occurs between bin $k = 58$ to $k = 62$. The peak of the histogram gives number of the pixels of the grass in the image. We call this number as x_g . From this, we compute the dominant grass pixel ratio (DGPR) as x_g/x , where x is the total number of pixels in the frame. We observed DGPR values vary from 0.16 to 0.24 for the field view. For non-field view image shown in Figure 3, we observed peak for bin $k = 20$ and its value is 2915. The image can be classified as field view or non-field view by using following condition:

```

if (DGPR > 0.16),
then frame belongs to class field view
else frame belongs to class non-field view

```

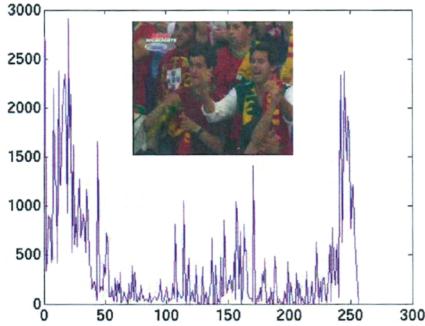


Fig. 3. Histogram of hue component of non-field view image (crowd), which is shown in the figure

C. Level-3a: Field View Classification

In order to reliably detect only the moving structures, the concept of the flux tensor is proposed in [17]. Flux tensor incorporates information about temporal gradient changes which leads to efficient discrimination between stationary and moving image features. Thus, the trace of the flux tensor matrix can be directly used to classify the moving and non-moving regions in the frames. Motion-mask is obtained by thresholding and post-processing averaged flux tensor trace. We used this motion-mask to classify the frame of the field view as long view, straight view and corner view. Our approach is summarized as follows:

Step-1: Generate motion-mask for the input image as shown in the second column of the Figure 4.

Step-2: Apply connected component technique to remove noisy objects from the image as shown in the third column of figure 4.

Step-3: In the connected component image, background color is the color of object 'field'. We divide the image into three parts 11, 12, and 2 as shown in the figure 4(a). Depending on the percentage of the field pixels (FP), we classify the frame into three classes using thresholds T_1 , T_2 and T_3 as follow:

```

if ( $FP_2 > T_1$ )  $\wedge$  (( $FP_{11} + FP_{12}$ )  $> T_2$ ),
then frame belongs to class long-view
else if  $|FP_{11} - FP_{12}| > T_3$ 
frame belongs to class corner-view
else
frame belongs to class straight-view

```

D. Level-3b: Close-up detection

Color-based object recognition is possible for sports video, since colors are purposely used to differentiate players of different teams, and umpire/referee. As shown in Figure 5, we divided the close-up frames into 16 blocks and observed that the skin pixels are more likely to occur in the block number 6, 7, 10 and 11. Depending upon, the skin percentage in these blocks, we can classify the frame as close-up or crowd. Our approach is summarized as follows:

Step-1: Convert the input $R - G - B$ image into $Y - C_b - C_r$

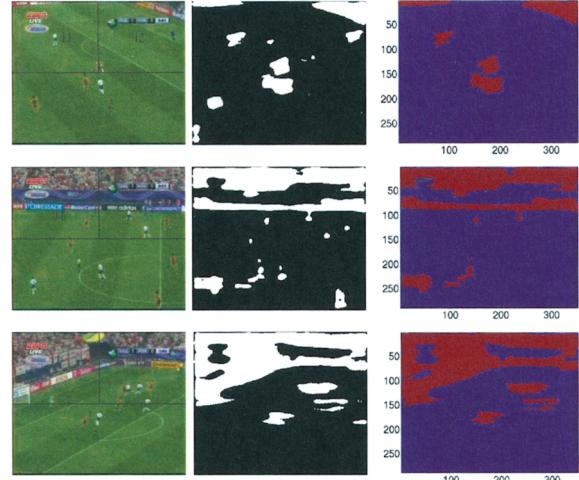


Fig. 4. Row-1 shows long view: (a) Image (b) motion-mask (c) connected component image, Row-2 shows straight view: (d) Image (e) motion-mask (f) connected component image, Row-3 shows corner view: (g) Image (h) motion-mask (i) connected component image

model. Use the following condition for detecting skin pixels.

```

if (105 <  $Y$  < 117)  $\wedge$  (110 <  $C_b$  < 113)  $\wedge$  ( $C_r$  > 128),
then pixel belongs to skin color
else pixel does not belong to skin color

```

Step-2: Apply connected component technique to remove noisy skin detected objects from the image as shown in Figure 6.

Step-3: Divide the image into 16 blocks, and compute the percentage of skin color pixels in each block.

Step-4: Let S_6 , S_7 , S_{10} , S_{11} are the skin percentages of block number 6, 7, 10, 11 respectively. Select threshold T_4 for considering the block as skin block. Apply the following condition to classify the image.

```

if ( $S_6 > T_4$ )  $\vee$  ( $S_7 > T_4$ )  $\vee$  ( $S_{10} > T_4$ )  $\vee$  ( $S_{11} > T_4$ ),
then frame belongs to class close up
else frame belongs to class crowd

```

E. Level-4: Close up Classification

In this level, we will classify the close-up images into player of team-A, player of team-B, goalkeeper of team-A, goalkeeper of team-B. We will first find out the location of the face of the person in the close-up using the approach used in level-3b. As shown in Figure 5, the face position will occur in the block 6,7,10,11. Depending on the face location, we will select the block for checking the jersey color of the player. For selecting the location of the jersey color, we have proposed the flow-chart as shown in Figure 7. Our approach is summarized as follows:

Step-1: Select the block for jersey color computation for the frame t using flowchart shown in Figure 7.

Step-2: Compute the mean of the red ($\mu_{r,t}$), mean of the green ($\mu_{g,t}$), mean of the blue ($\mu_{b,t}$) components of the pixels in that block of frame t .

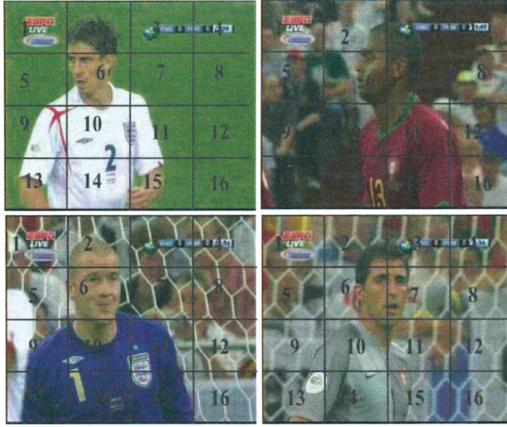


Fig. 5. Close-up frames frequently observed in broadcasted football video. Row-1: (a)Player of Team-A, (b) Player of Team-B, Row-2: (c) Goalkeeper of Team-A, (d) Goalkeeper of Team-B

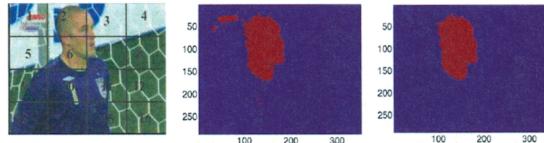


Fig. 6. (a) Image of goalkeeper (b)Image showing skin detection, (c) Connected component image

Step-3: Compute the mean of the red ($\mu_{r,k}$), mean of the green ($\mu_{g,k}$), mean of the blue ($\mu_{b,k}$) components of the pixels of the block of the known class k . Let C be the total number of classes. In our case, the four classes ($C = 4$) are: player of team-A (class-1), player of team-B (class-2), goalkeeper of team-A (class-3), goalkeeper of team-B (class-4).

Step-4: Compute the Euclidean distance of the block of frame t from the class k using following formula.

$$D_{k,t} = \sqrt{(\mu_{r,t} - \mu_{r,k})^2 + (\mu_{g,t} - \mu_{g,k})^2 + (\mu_{b,t} - \mu_{b,k})^2} \quad (5)$$

for $k = 1$ to C

Step-5: Select the value of k for which $D_{k,t}$ has lowest value. Assign that value of k as a class-label to the particular frame t .

F. Event

We have defined the events as scenes in the video with some semantic meaning (i.e. labels from a semantic hierarchy) attached to it based on the leaf nodes shown in Figure 1. Events are extracted as the leaf nodes of the level-3 and 4 of hierarchical tree. The events are long view, straight view, corner view, crowd, player of team-A, player of team-B, goalkeeper of team-A, goalkeeper of team-B.

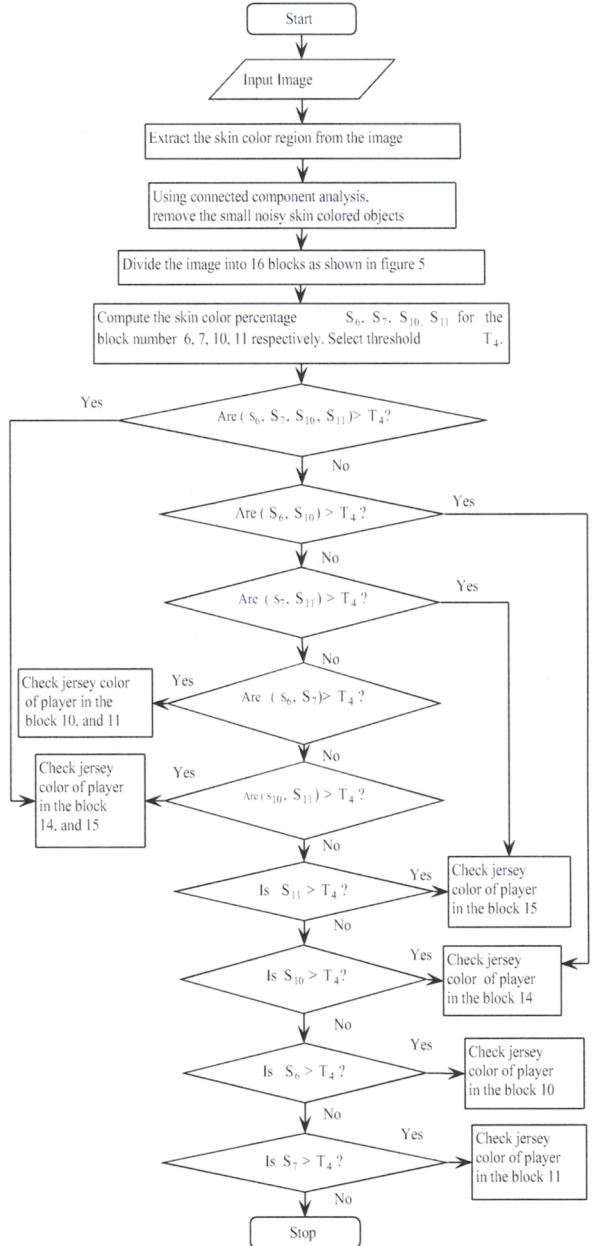


Fig. 7. Flow chart for the block selection for checking player's jersey color

TABLE I
SOCCER VIDEO SEQUENCES USED FOR TESTING

ID	Name of the match	Broad -cast channel	A vs B	Date	Match Result
V_1	FIFA-world Cup Quarter final	ESPN	Germany vs Argentina	30/06/ 2006	Germany won
V_2	FIFA-world Cup Quarter final	ESPN	England vs Portugal	01/07/ 2006	Portugal won
V_3	FIFA-world Cup Quarter final	ESPN	Brazil vs France	01/07/ 2006	France won

TABLE II
PERFORMANCE OF LEVEL-1 OF HIERARCHICAL CLASSIFIER

Video ID	Total Duration	Extracted Clips Duration	Recall	Precision
V_1	90 min	54 min	87.68 %	73.76 %
V_2	90 min	52 min	100 %	93 %
V_3	90 min	58 min	100 %	81 %

III. EXPERIMENTAL RESULTS

We have tested our proposed approach using live recordings of three FIFA world cup matches as shown in table I. Since commercials may also be classified as a excitement based on audio excitement, we remove the commercials from the sports video before applying hierarchical classifier. We define the threshold vector as $T = [T_1, T_2, T_3, T_4]$. We experimentally found that $T = [65, 65, 10, 60]$ gives better results. For level-1, we are classifying the segments of the video, and level-2 to level-4, we are classifying the frames of the excited segment. For measuring the performance of classifiers at each level, we use following parameters:

$$Recall = \frac{N_c}{N_c + N_m} \quad (6)$$

$$Precision = \frac{N_c}{N_c + N_f} \quad (7)$$

Where, N_c , N_m , N_f represents the number of excitement clips correctly detected, missed and false positive, respectively, for level-1. For level-2 to level-4, N_c , N_m , N_f represents the number of frames correctly detected, missed and false positive, respectively.

A. Level-1:Excitement Detection

We present here the results of typical soccer video clip of duration 2 minutes 25 seconds. We sampled audio at a rate of 44.1 KHz and video frames are down sampled from its original 30 frames/second to 5 frames/second. We extracted 725 video frames of this video clip. We selected excitement clip from video frame number 131 to 288 on the basis of the product of short-time audio energy and ZCR as shown in Fig 8(d).

The overall performance of the classifier at level-1 is shown in table II. In case of poor broadcasting quality and noisy audio, performance of audio-based excitement clip extraction decreases.

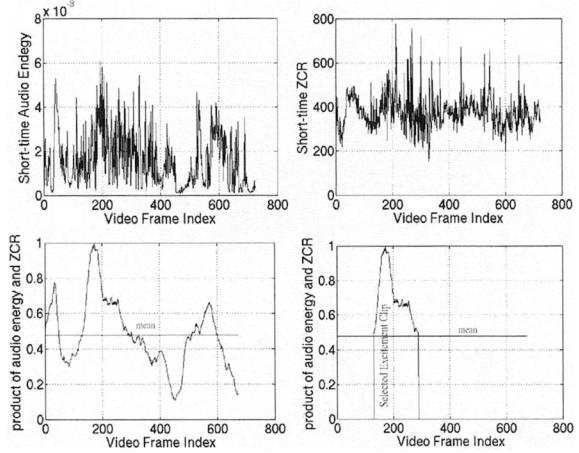


Fig. 8. (a)Audio energy vs Video Frame Number (b) ZCR vs Video Frame Number (c) Product of Audio energy and ZCR vs Video Frame Number (d) Frames with product greater than Threshold and of duration greater than 100 frames are selected



Fig. 9. Row-1: (a) Field view (b) Field view (c) Field view, Row-2: (d) Non-field view (e) Non-field view (f) Non-field view

B. Level-2:Field View Detection

Table III shows the classification performance of field view detection of the images shown in Figure 9. The $DGPR$ values for the grass lies between 0.16 to 0.24. For non-field view, we observed the values less than 0.1. Because of this discrimination, we observed almost 100 % recall and precision.

C. Level-3a:Field View Classification

The overall performance of the classifier at level-3a is shown in table IV. Since we consider few previous frame for generating motion-mask, we observed some missclassification near the shot boundaries.

D. Level-3b:Close-up Detection

Table V shows the performance of close-up detection for the images of Figure 5. The performance of this classifier is tested with 2005 frames as shown in the table VI.

TABLE III
PERFORMANCE OF LEVEL-2 OF HIERARCHICAL CLASSIFIER

Image	DGPR	Actual class	Observed Class
Figure 9(a)	0.3206	field view	field view
Figure 9(b)	0.2196	field view	field view
Figure 9(c)	0.2254	field view	field view
Figure 9(d)	0.0028	non-field view	non-field view
Figure 9(e)	0.0015	non-field view	non-field view
Figure 9(f)	0.0031	non-field view	non-field view

TABLE IV
PERFORMANCE OF FIELD-VIEW CLASSIFICATION AT LEVEL-3A

Total Frames	N _c	N _m	N _f	Recall	Precision
2125	2031	94	102	95.58 %	95.22 %

E. Level-4: Close-up Classification

We have trained our classifier for classifying the close-up frames into four classes. The classes are player of team-A, player of team-B, goalkeeper of team-A, goalkeeper of team-B. We tested 1375 close-up frames and observed the classification performance as shown in table VII. We observed less precision, because the close-up of the person from spectator gets classified into one of the classes.

IV. CONCLUSION

In this paper, we have presented a hierarchical framework for analyzing high-level events in soccer video by combining low level feature analysis with high level semantic knowledge. The sports domain semantic knowledge encoded in the hierarchical classification not only reduces the cost of processing data drastically, but also significantly increases the classifier accuracy. The hierarchical framework enables the use of simple features and organizes the set of features in a semantically meaningful way. The proposed hierarchical semantic framework for event classification can be readily generalized to other sports domains as well as other types of video. Our future work includes probabilistic modeling framework for building the semantic hierarchy, semantic concept extraction based on the classified scenes (events) for highlight generation and video summarization.

TABLE V
CLASSIFICATION OF NON-FIELD VIEWS INTO CROWD AND CLOSE-UPS AT LEVEL-3B

Image	S ₆	S ₇	S ₁₀	S ₁₁	Actual	Observed
Fig. 5(a)	61	0	70	0	close-up	close-up
Fig. 5(b)	71	91	62	94	close-up	close-up
Fig. 5(c)	85	43	12	2	close-up	close-up
Fig. 5(d)	58	73	17	59	close-up	close-up

TABLE VI
PERFORMANCE OF CLOSE-UP DETECTION AT LEVEL-3B CLASSIFIER

Total Frames	N _c	N _m	N _f	Recall	Precision
2005	1807	198	203	90.12%	89.90%

TABLE VII
PERFORMANCE OF CLOSE-UP CLASSIFICATION AT LEVEL-4

Total Frames	N _c	N _m	N _f	Recall	Precision
1375	1196	179	192	86.98%	85.8%

REFERENCES

- [1] A. Kokaram, N. Rea, R. Dahyot, M. Tekalp, P. Bouthemy, P. Gros and I. Sezan, *Browsing sports video: trends in sports-related indexing and retrieval work*, in IEEE Signal Processing Magazine, vol. 23, no. 2, pp. 47-58, 2006.
- [2] Y. Li, J. Smith, T. Zhang and S. Chang, *Multimedia Database Management Systems*, in Elsevier Journal of Visual Communication and Image Representation, pp. 261-264, 2004.
- [3] J. Wang, E. Chng, C. Xu, H. Lu and Q. Tian, *Generation of Personalized Music Sports Video Using Multimodal Cues*, in IEEE Transaction on Multimedia, vol. 9, no. 3, pp. 576-588, 2007.
- [4] G. Zhu, Q. Huang, C. Xu, L. Xing, W. Gao and H. Yao, *Human Behavior Analysis for Highlight Ranking in Broadcast Racket Sports Video*, in IEEE Transactions on Multimedia, vol. 9, no. 6, pp. 1167-1182, 2007.
- [5] M.H. Kolekar and S. Sengupta, *Event-importance Based Customized and Automatic Cricket Highlight Generation*, IEEE Int. Conf. on Multimedia and Expo, 2006.
- [6] C. Xu, J. Wang, H. Lu and Y. Zhang, *A Novel Framework for Semantic Annotation and Personalized Retrieval of Sports Video*, in IEEE Transactions on Multimedia, vol. 10, no. 3, pp. 421-436, 2008.
- [7] L. Duan, M. Xu, Q. Tian, C. Xu and J. Jin, *A Unified Framework for Semantic Shot Classification in Sports Video*, in IEEE Transactions on Multimedia, vol. 7, no. 6, pp. 1066-1083, 2005.
- [8] A. Hanjalic, *Generic Approach to Highlight Extraction from a Sport Video*, in proc. ICIP, vol. 1, pp. 1-4, 2003.
- [9] B. Li, H. Pan and I. Sezan, *A General Framework for Sports Video Summarization with its application to Soccer*, in IEEE Int. Conf. on Audio, Speech, and Signal Processing, pp. 169-172, 2002.
- [10] N. Babaguchi, Y. Kawai, T. Ogura and T. Kitahashi, *Personalized Abstraction of Broadcasted American Football Video by Highlight Selection*, in IEEE Transaction Multimedia, vol. 6, no. 4, pp. 107-109, 2004.
- [11] K. Wan and C. Xu, *Recent soccer highlight generation with a novel dominant speech feature extractor*, in IEEE Int. Conf. on Multimedia and Expo, vol. 1, pp. 591-594, 2004.
- [12] Y. Ding and G. Fan, *Segmental Hidden Markov Model for View-based Sports Video Analysis*, in IEEE Conf. on Computer Vision and Pattern Recognition, pp. 17-22, 2007.
- [13] V. Chasanis, A. Likas, and N. Galatsanos, *Scene Detection in Videos Using Shot Clustering and Symbolic Sequence Segmentation*, in IEEE Workshop on Multimedia Signal Processing, pp. 187-190, 2007.
- [14] C. Ngo, T. Pong, H. Zhang, *On clustering and retrieval of video shots through temporal slices analysis*, IEEE Transactions on Multimedia, vol. 4, no. 4, pp. 446-458, 2002.
- [15] M. H. Kolekar and S. Sengupta, *Semantic Indexing of News Video Sequences: A Multimodal Hierarchical Approach Based on Hidden Markov Model*, in Proc. of IEEE International Region 10 Conference (TENCON), Melbourne, Australia, Nov 2005
- [16] P. Xu, L. Xie, S. Chang, A. Divakaran, A. Vetro, H. Sun, *Algorithms and system for segmentation and structure analysis in soccer video*, IEEE Int. Conf. on Multimedia and Expo, pp. 721-724, 2001.
- [17] F. Bunyak, K. Palaniappan, S. Nath, and G. Seetharaman, *Flux Tensor Constrained Geodesic Active Contours with Sensor Fusion for Persistent Object Tracking*, in Journal of Multimedia, vol. 2, no. 4, pp. 20-33, 2007.