

Shot Boundary Detection and Keyframe Extraction based on Scale Invariant Feature Transform

Gentao Liu

Xiangming Wen

Wei Zheng

Peizhou He

School of Information and Communication Engineering,
Beijing University of Posts and Telecommunications,
Beijing, China

liugentao@gmail.com

xiangmw@bupt.edu.cn

zhengweius@gmail.com

hepeizhou@gmail.com

Abstract—In shot boundary detection, the key technology is to compute the visual content discontinuity values between consecutive video frames. In this paper, a unified framework is proposed to detect the shot boundaries and extract the keyframes of a shot. Firstly, the Scale invariant feature transform (SIFT) is adopted to compute the visual content discontinuity values. Then a new method, which is called the Local Double Threshold Shot Boundary Detection (LDT-SBD), is used to detect shot boundaries. Lastly, two mechanisms are proposed to extract keyframe. Experimental results show the framework is effective and has a good performance.

Keywords- Shot boundary detection; Keyframe extraction; Scale Invariant Feature Transform; LDT-SBD

I. INTRODUCTION

With the increasing proliferation of digital videos in recent years, efficient techniques based on video contents have become the major methods in the indexing and retrieval systems of video. Consequently, the video content analysis has been a very active research field.

Parsing a video into its basic temporal units – shots – is considered as the initial step in the process of video content analysis. A shot can be defined as a sequence of frames generated during a continuous camera operation and represents a continuous action in time and space [1]. Two consecutive shots are separated by a shot boundary. There are two categories of shot boundaries: cut and gradual. Cut is generated by simply attaching one shot to another without modifying them while a gradual shot boundary is the result of apply an editing effect to merge two shots, which includes dissolves, fades, wipe and etc. Over the last decades, shot boundary detection has received a considerable amount of research attention and a large number of algorithms have been developed [2-4]. However, most of the algorithms are sensitive to noise and movements of cameras or targets and their performance are still insufficient due to lacking of image content information. Some researchers are making effort to find techniques that can improve the performance of shot boundary detection, such as using scale invariant feature transform (SIFT).

SIFT was first presented by David G.Lowe in 2004 and has been well developed for object recognition and matching because its image features are invariant to image rotation, scale and robust across a substantial range of affine

distortion, addition of noise, and change in illumination [5]. As in the field of video parsing, it's mainly used for motion estimation and object detection. The SIFT is applied to shot boundary detection by Min-Ho Park [6]. He uses the number of matched SIFT features between two adjacent frames to detect cut shot boundary and then the matched number between two frames with a distance of 22 is used to detect gradual shot boundary. The method attains a satisfactory performance in cut detection, but lacks robust in detecting gradual shot boundary because of the fixed frame distance. In addition, this method detects cut and gradual in two steps and needs more computing and time cost.

After the video is segmented into shots, the most representative frames in each shot, named keyframe, are extracted for further applications, such as video retrieval, user browsing and content analysis. So the extraction method of keyframe in a shot is another key study direction. The easiest method is to take the first, middle or ending frame as keyframe. This method results in one possibility that the extracted keyframes have a low correlation in visual content. Another way is to split a shot into smaller video clips. Each clip can be a fixed time video or consecutive frame series characterized by high vision-content redundancy. Each clip typically includes several consecutive frames, but may also stretch to a complete shot. Then we can extract keyframes from the clips in the way that the visual content redundancy is minimized [7].

As the above mentioned analysis, this paper presents a unified framework for shot boundary detection and keyframe extraction. Firstly, ratio of matched feature number to their total number is computed to measure the similarity between two adjacent frames. Then a new method named Local Double Threshold Shot Boundary Detection (LDT-SBD) is proposed to detect shot boundaries. This method can avoid false detection caused by few SIFT keypoints generated and is robust to noise, camera rotation and varying length of the gradual transition and the cut and gradual shot boundaries can be detected in the same time. For the purpose of keyframe extraction, we put forward two mechanisms to split shot into clips according to content coherence. Chain-matched keypoints number among consecutive frames indicates the visual content coherence of these frames. If the number is sufficient low, we could suppose that the content has changed, a new clip boundary found, and the chain-matching procedure restarted. Another mechanism is to use

the pattern of matching ratio change. Once the high content coherent clip is produced, the frame with most SIFT keypoints can be extracted as a keyframe since more keypoints means more vision contents and the visual content redundancy is minimized.

The reminder of this paper is organized as follows. In Section II, the proposed framework is described in detail. Experimental results and discussions are shown in Section III. Finally conclusions are given in Section IV.

II. THE PROPOSED FRAMEWORK

A. Framework

The framework of shot boundary detection and keyframe extraction using Scale Invariant Feature Transform can be illustrated as Fig. 1.

SIFT keypoints are extracted from each frame of videos in terms of their temporal Sequence. Each keypoint is represented by a descriptor of 128 dimensions. Then a Best-Bin-First (BBF) [5] algorithm is used to match keypoints between two adjacent frames. And then ratios of match keypoints number to total number are used to detect shot boundaries and extract keyframes.

B. SIFT

SIFT is an approach for detecting and extracting distinctive invariant feature descriptors from image. There are four major steps: detection of scale-space extreme, accurate keypoint localization, orientation assignment, descriptor representation.

1) *Detection of scale-space extreme* [8]: Key point candidates for SIFT features are obtained potentially from local extrema of difference-of-Gaussian (DoG) space. DoG scale space can be obtained from (1).

$$D(x, y, \sigma) = (G(x, y, k\sigma) - G(x, y, \sigma)) * I(x, y) \quad (1)$$

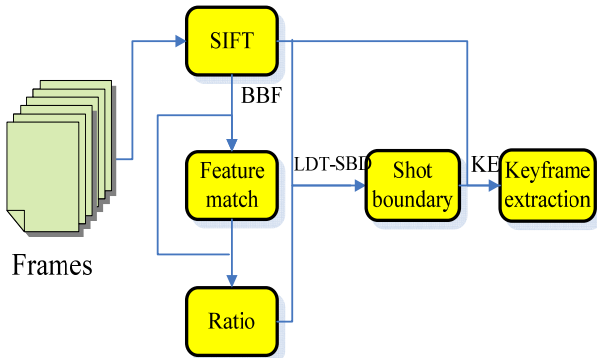


Figure 1. The framework of shot boundary detection and keyframe extraction

where $*$ is the convolution operation, and $I(x, y)$ is the gray value of pixel at (x, y) . $G(x, y, \sigma)$ is a variable-scale Gaussian kernel defined as:

$$G(x, y, \sigma) = \frac{1}{2\pi\sigma^2} e^{-(x^2+y^2)/2\sigma^2} \quad (2)$$

2) *Accurate keypoint localization*: Candidate keypoints are localized to sub-pixel accuracy and low contrast keypoints introduced by noise and edge response will be removed.

3) *Orientation assignment*: Each keypoint is assigned a orientation that makes sure the descriptor invariant to rotation. The orientation can be determined by computing an orientation histogram from the gradient orientation of sample points within a region around the keypoint.

4) *Descriptor representation*: The feature vector descriptor is computed as a set of orientation histograms in the 16×16 pixel region. Since every 4×4 region is projected onto one histogram and each histogram contains 8 bins for 8 directions, it leads to a SIFT feature vector of 128 dimensions. The vector is normalized to enhance invariance to changes in illumination.

After the SIFT feature descriptors are generated, we can measure the similarity of two descriptor vectors by the Euclidean distance. Two frame's SIFT features matching can be implemented with a BBF algorithm.

C. Shot boundary detection

Our work on detecting shot boundaries is based on the basic assumption that frames surrounding a boundary are, in general, much more different with respect to their visual content than the frames taken from within a shot [7]. So the problem of detecting shot boundary may be approached by searching for large discontinuities in the flow of videos and selecting suitable features to quantify the discontinuity. Digital videos generally have a high frame rate of more than 20 frames per second; therefore, two consecutive frames of one shot are likely to contain a considerable portion of the same material. This can be reflected by the matched SIFT keypoints between the very two frames.

Reference [6] uses the number of matched point to depict the common material between consecutive frames and detect shot boundary. The number relates to not only common material but also total number of SIFT feature keypoints. If frames looks simple, or with few colors, we would get fewer matched keypoints even though they are similar. Therefore, we use the ratio of matched number to total number instead to avoid false detection caused by too few keypoints generated and it is defined as:

$$R(t) = \frac{2M(t)}{F(t-1) + F(t)} \quad (3)$$

where $F(t)$ and $F(t-1)$ are the amount of keypoints generated from frame t and $t-1$. $M(t)$ is the matched keypoints number between these two frames. A $R(t)$ -curve is generated according to a video clip of “Tom and Jerry” and presented in Fig. 2. Dash-line represents cut while dot-line represents gradual shot boundary.

We propose a new method called Local Double Threshold Shot Boundary Detection (LDT-SBD) using the $R(t)$ curve to detect shot boundary in a unify frame which can be carried out as follows:

1) *Moving average calculation*: Calculate the moving average value of frame t as:

$$\overline{R(t)} = \frac{1}{K} \sum_{i=t-K}^{t-1} R(i) \quad (4)$$

where K is the length of frames used to calculate the moving average value. $\overline{R(t)}$ depicts the local average of $R(t)$.

2) *Detect shot boundary*: Frame content changes drastically at shot boundary and distinguishable local minima of the $R(t)$ curve indicates the very changes. We can measure the change with the difference of $R(t)$ and $\overline{R(t)}$. If

$$\overline{R(t)} - R(t) \geq \Delta \quad (5)$$

then the frame t is a cut shot boundary. Else if

$$\delta \leq \overline{R(t)} - R(t) < \Delta \quad (6)$$

at several consecutive frames, there may be a gradual shot boundary. Further the sum of these differences can be computed. If

$$\sum_{t=i}^{i+j} (\overline{R(t)} - R(t)) \geq \Delta' \quad (7)$$

then it is a gradual shot boundary.

D. Keyframe extraction

The keyframe set which is to be extracted for the purpose of shot representation needs to satisfy the following two requirements [7]:

- The redundancy in the vision content captured by the key frames is minimized.
- Keyframes capture all relevant aspects of the visual content

Obviously, the amount of keyframes extracted to represent a shot and the position of the keyframe depend on the vision content of the shot. As mentioned above, this vision content can be reflected by SIFT keypoints. We

propose two mechanisms to segment shot into smaller temporal units-clip according to vision content firstly.

1) *Chain-matching*: SIFT matching process is carried out sequentially between features of a new frame and matched features prior. As the matching proceeds, the amount of matched keypoints will decrease to 0, a clip boundary is found. The matched keypoints number indicates the common feature and content similarity among frames.

2) *Rule-matching*: We can find out that the $R(t)$ -curve has a corresponding relationship with chain-matching result from the Fig. 3, where solid-line is the $R(t)$ -curve, dot-line is the chain-matching result, dash-lines represent shot boundaries and dash-dot-lines represent clip boundaries. We use two rules to describe this relationship.

- Rule 1: Calculate the local average of $R(t)$. If there is a sharp variation at frame k , it is determined as a clip boundary.
- Rule 2: Calculate the secondary moment of $R(t)$ and sharp variation means clip boundary.

These clips are high content redundancy and one frame is sufficient to represent an entire clip. We select the frame with most keypoints as the key frame since more keypoints

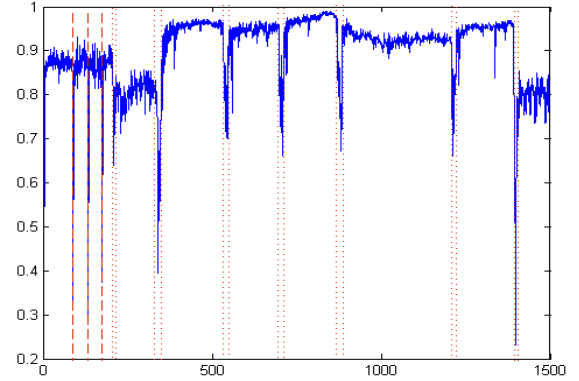


Figure 2. A sample of $R(t)$ -curve

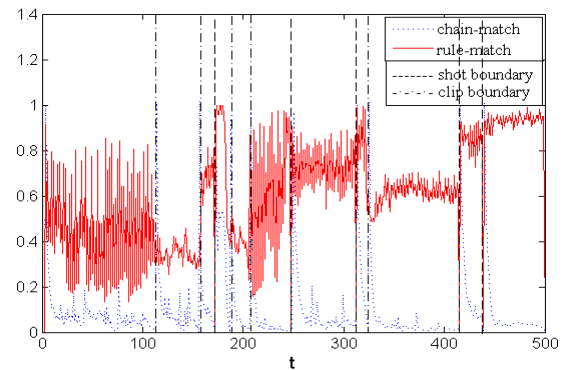


Figure 3. Chain-matching result

means more vision contents and the frame is most representative.

III. EXPERIMENTAL RESULTS

We have designed an experimental system for evaluating the performance of our method. Four types of videos are used for shot boundary detection and keyframe detection, including news, movie, sports and cartoon. The shot boundaries that are used as the ground truth are labeled manually at the beginning of the experimental. The information of these experimental materials is listed in table I.

TABLE I. NUMBERS OF FRAMES AND SHOT BOUNDARIES OF EXPERIMENTAL MATERIALS

Video type	Number of frames	Cut	Gradual
News	5929	21	13
Movie	6320	51	14
Sports	5995	40	7
Cartoon	3936	37	2

A. Shot boundary detection

We use precision, recall, which are defined as (8) and (9), to measure the performance of shot boundary detection.

$$Precision = \frac{N_d}{N_d + N_f} \quad (8)$$

$$Recall = \frac{N_d}{N_d + N_m} \quad (9)$$

where N_d , N_f and N_m are the numbers of correct, false and miss shot boundary detections, respectively. The experimental result of proposed LDT-SBD algorithm is shown in Table II and it shows that the method is effective for all of the four types of video.

TABLE II. PERFORMANCE OF LDT-SBD ALGORITHM

Video type	Cut		Gradual		Average	
	Precision	Recall	Precision	Recall	Precision	Recall
News	0.952	0.952	1.000	0.923	0.969	0.941
Movie	0.960	0.941	0.813	0.929	0.924	0.938
Sports	0.946	0.875	0.750	0.857	0.911	0.872
Cartoon	0.923	0.972	0.667	1.000	0.905	0.974

We also have compared our method to the algorithm proposed in [6] and the results are showed in table III. From the table III, we can see that our method get an obvious improvement in detecting gradual shot boundary.

B. Keyframe extraction

We perform our experiments using the two mechanisms proposed above and in most case, they get a similar result.



Figure 4. Keyframes extracted from “Tom and Jerry”

TABLE III. COMPARISON BETWEEN TWO METHODS

	Cut		Gradual		Average	
	Precision	Recall	Precision	Recall	Precision	Recall
Our method	0.946	0.933	0.868	0.916	0.930	0.930
Method in [6]	0.946	0.933	0.783	0.805	0.913	0.908

Fig.4 presents a sequence of shot keyframes we extract from the video “Tom and Jerry” according to Fig. 3. The 1st, 2nd and 3rd keyframe are extracted from Shot 1; the 4th, 5th, 6th are extracted from Shot 2; 8th, 9th are extracted from Shot 4. As the vision contents of Shot 3, 5, 6 have less change, only one frame is extracted.

IV. CONCLUSION

In this paper, we propose a novel framework for shot boundary detection and keyframe extraction base on SIFT. Firstly we detect two types of shot boundaries-cut and gradual-using a LDT-SBD algorithm. Then the shots segmented to smaller temporal clip according to their content similarity and the most representative frame in a clip is extracted as keyframe. Experimental results show that the framework gives satisfactory performance for both shot boundary detection and keyframe extraction.

Future research will focus on the reduction of miss or false shot boundary detection and investigation this method for video index and content analysis.

V. ACKNOWLEDGMENT

This work was supported by National Natural Science Foundation of China (60743007, 60872050) and Beijing Municipal Education Commission under a Special Grants for Building Projects (XK100130648).

REFERENCES

- [1] G. Davenport, T. A. Smith, and N. Pinciver, "Cinematic primitives for multimedia," IEEE Computer Graphics and Applications, vol. 11, pp. 67-74, 1991.
- [2] X. Ling, L. Chao, L. Huan, and X. Zhang, "A general method for shot boundary detection," Busan, South Korea, 2008, pp. 394-397.

- [3] S.-T. Chiu, G.-S. Lin, M.-K. Chang, and H.-Y. Wu, "An effective shot boundary detection algorithm for movies and sports," Dalian, Liaoning, China, 2008, pp. 460-462.
- [4] S. Lefevre and N. Vincent, "Efficient and robust shot change detection," *Journal of Real-Time Image Processing*, vol. 2, pp. 23-34, 2007.
- [5] D. G. Lowe, "Distinctive image features from scale-invariant keypoints," *International Journal of Computer Vision*, vol. 60, pp. 91-110, 2004.
- [6] M.-H. Park, R.-H. Park, and S. W. Lee, "Shot boundary detection using scale Invariant feature matching," San Jose, CA, United States, 2006, p. 60771.
- [7] A. Hanjalic, *content-based analysis of digital video*: Kluwer Academic Publishers 2004.
- [8] A. P. Witkin, "SCALE-SPACE FILTERING: A NEW APPROACH TO MULTI-SCALE DESCRIPTION," San Diego, CA, USA, 1984, pp. 39-41.