

# A Review of Real-Time Segmentation of Uncompressed Video Sequences for Content-Based Search and Retrieval

Sébastien Lefèvre, Jérôme Holler, Nicole Vincent

► **To cite this version:**

Sébastien Lefèvre, Jérôme Holler, Nicole Vincent. A Review of Real-Time Segmentation of Uncompressed Video Sequences for Content-Based Search and Retrieval. *Real-Time Imaging*, 2003, 9 (1), pp.73-98. <10.1016/S1077-2014(02)00115-8>. <hal-00512620>

**HAL Id: hal-00512620**

**<https://hal.archives-ouvertes.fr/hal-00512620>**

Submitted on 31 Aug 2010

**HAL** is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

# A Review of Real-time Segmentation of Uncompressed Video Sequences for Content-Based Search and Retrieval\*

Sébastien Lefèvre<sup>1,2</sup>, Jérôme Holler<sup>1</sup>, Nicole Vincent<sup>1</sup>

<sup>1</sup> Laboratoire d'Informatique  
E3i / Université de Tours  
64 avenue Portalis, 37200 Tours - FRANCE

<sup>2</sup> AtosOrigin  
19 rue de la Vallée Maillard  
BP1311, 41013 Blois Cedex - FRANCE  
(lefevre,vincent)@univ-tours.fr

## Abstract

We present in this paper a review of methods for segmentation of uncompressed video sequences. Video segmentation is usually performed in the temporal domain by shot change detection. In case of real-time segmentation, computational complexity is one of the criteria which has to be taken into account when comparing different methods. When dealing with uncompressed video sequences, this criterion is even more significant. However previous published reviews did not involve complexity criterion when comparing shot change detection methods. Only recognition rate and ability to classify detected shot changes were considered. So contrary to previous reviews we give here complexity of most of the described methods. We review in this paper an extensive set of methods presented in the literature and classify them in several parts, depending on the information used to detect shot changes. The earliest methods were comparing successive frames by relying on most simple elements, that is to say pixels. Comparison could be performed on a global level, so methods based on histograms were also proposed. Block-based methods have been considered to process data at an intermediate level, between local (using pixels) and global (using histograms) levels. More complex features can be involved, resulting in feature-based methods. Alternatively some methods rely on motion as a criterion to detect shot changes. Finally different kinds of information could be combined together in order to increase the quality of shot change detection. So our review will detail segmentation methods

---

\*S. Lefèvre, J. Holler, N. Vincent, A Review of Real-Time Segmentation of Uncompressed Video Sequences for Content-Based Search and Retrieval, Real Time Imaging, Vol. 9, No. 1, 2003, pages 73-98, doi:10.1016/S1077-2014(02)00115-8

based on the following information: pixel, histogram, block, feature, motion, or other kind of information.

## 1 Introduction

Multimedia information is more and more used, mostly thanks to increasing computation resources. One of the main processings needed when dealing with multimedia data is multimedia sequence indexing. The importance of this research field can be shown by the number of recent communications and publications on the subject. In order to index multimedia data, we may need a preprocessing the aim of which is to temporally segment the videos, that is to say detect the shot changes present in the video sequences.

The number of shot change detection methods is now important and several reviews of these methods have been done [3, 42, 17, 7, 48, 57, 44, 14, 56, 22, 88, 31, 32, 59]. These reviews often present the different methods and their efficiency based on some quality measures. So they are very useful when one wants to select and implement a shot change detection method for a global video processing which could be done off-line. When processing has to be done on-line, the selection of a particular method should also consider computation time. This is especially true when dealing with uncompressed video sequences which contain a huge quantity of data. If the method has to be implemented on common hardware architecture, computation time is directly linked with complexity of the method. So in this paper we are to review most of the methods presented in the literature and to focus on their complexity.

In the first part of this paper, before we present a large number of methods, we are to situate our contribution. We are giving some references to previous reviews on the subject. It is also necessary to recall and describe the different forms a shot change can take. We give a few details on quality evaluation and introduce the way we compute the complexity of the methods. We also define the notations used in this paper. The next parts will deal with a description of the encountered methods. Finally some conclusions will be given about shot change detection methods.

### 1.1 Related works

Several reviews have already been published in the literature. Ahanger et al [3] discuss requirements and global architectures for video indexing. Some video segmentation methods are presented in this framework. Idris et al [42] are dealing with image and video indexation. Image features and video processing algorithms useful for indexation are described. Shot change detection is one of the video processing needed to characterize a video sequence. Brunelli et al [17] also present video indexation, and describe main algorithms including shot change detection. They are particularly involved in a video indexing system. Aigrain et al [7] review techniques for video content analysis. Shot change detection is one of these techniques. Koprinska et al [48] review algorithms for shot change detection and camera operation recognition. Lienhart [57] presents the different kinds of shot changes and some dedicated detection methods.

Jiang et al [44] propose a review based on three categories which are uncompressed video-based, compressed video-based, and model-based algorithms.

Some reviews are comparing a few algorithms based on author's implementation. Boreczky et al [14] compare the performances of 5 algorithms using a common evaluation methodology. Lienhart [56] compares the methods and characterizes their ability to correctly determine the kind of shot changes that have been detected. Dailianas et al [22] compare several segmentation methods and introduce a filtering algorithm in order to limit the false detections. Some information is also given about complexity of the evaluated methods. Yusoff et al [88] compare several methods and propose improved versions using an adaptive threshold.

Finally some papers are reviewing only a specific kind of methods, as those from Gargi et al [31, 32] which are respectively dedicated to colour histogram-based methods and MPEG and motion-based methods for temporal segmentation of video. Mandal et al [59] focused on methods working in compressed domain.

Contrary to other approaches, we review and compare in this paper uncompressed video segmentation methods following their computational complexity and not their detection or error rates, which has been already done in papers presented in this section. We base the classification of the presented methods on the basic elements used in the segmentation process: pixels, histograms, blocks, features, motion, and combination of several approaches.

## 1.2 Shot change description

A shot is defined as a continuous video acquisition (with the same camera). When the video acquisition is done with another camera, there is a shot change. The simplest way to perform a change between two shots is called a *cut*. In this case, the last frame of the first video sequence is directly followed by the first frame of the second video sequence. This kind of shot change is also called abrupt change. Because of their simplicity, cuts are often the easiest shot changes to be detected.

More complex shot changes are now available for video editing, thanks to improvement of the video production softwares. Instead of cutting and pasting the second video next to the first one, it is possible to insert an effect, as a *wipe*, a *fade*, or a *dissolve*. A wipe is obtained by progressively replacing the old image by the new one, using a spatial basis. A dissolve is a transition where all the images inserted between the two video sequences contain pixels whose values are computed as linear combination of the final frame of the first video sequence and the initial frame of the second video sequence. Fades are special cases of dissolve effects, where a monochrome frame replaces the last frame of the first shot (*fade in*) or the first frame of the second shot (*fade out*). There are also other kinds of effects (combining for example wipe and zoom), but actually most of the shot change detection methods are concerned only by the effects described previously in their indexing task.

## 1.3 Quality evaluation

The recognition rate is the most used quality criterion in order to compare shot change detection methods. Some work has been done to define some standard quality measures

and to discuss existing ones [31, 27, 37, 72]. Most of the time, quality is evaluated thanks to computation of the quantity of correctly detected shot changes, missed shot changes, and false detections. Indeed, to be fair, the evaluation should be achieved on a universal benchmark. This is not the case.

A similar background for two consecutive shots results often in missing the shot change. False detections appear when there is significant content change. Camera motion, moving objects, illumination changes can be sources of false detections.

Our review does not focus on this aspect of quality evaluation of the methods. We will insist here on complexity of the methods, because this work has not been done for a complete set of methods yet.

In order to compare video segmentation methods, it is also possible to take into account the number of thresholds or parameters which have to be set. Learning capabilities of these thresholds or parameters can also be used as comparison criteria. Comparison of uncompressed video segmentation methods based on these criteria is out of the scope of this paper.

## 1.4 Complexity computation

As mentioned previously, quality evaluation is not the only criterion to evaluate and compare shot change detection methods if we are concerned with real-time (or near real-time) processing using common hardware. In this case, one should also have to consider complexity of the evaluated methods. A work on complexity of shot change detection methods has been done by Dailianas et al in [22] but it was limited to few methods.

In this paper, the complexity was computed considering a cost of one for any logical or arithmetic operation (including absolute value). We do not consider other operations as for example memory access time or branching operations (*e.g.* if ... then). In order to compute the complexity of the methods, we define  $\mathcal{N}$  as the possible number of levels for pixel value, which is equivalent to the number of bins for histogram-based methods. We also introduced  $\mathcal{P}$  as the number of pixels per frame. In case of block-based methods, we use the notation  $\mathcal{B}$  to represent the number of blocks defined in the frame.

Complexity measurements given in this paper represent the number of operations needed to process one frame. Temporal subsampling of the video sequences is not taken into account. However, when values obtained for a given frame can be used to process the next frame, complexity measurements are optimised and given considering the use of previous results.

## 1.5 Notations

Video sequences are composed of successive frames or images. We define  $I_t$  the frame of the video obtained at time  $t$ . So it is possible to define  $P(I_t, i, j)$  the intensity of the pixel with coordinates  $i$  and  $j$  in the frame  $I_t$ . We assume that the size of the images is  $X$ -by- $Y$  pixels, so we have  $1 \leq i \leq X$  and  $1 \leq j \leq Y$ .

When methods are dealing with colour images, the notation  $P(I_t, C_k, i, j)$  will be used.  $C_k$  represents the colour component numbered  $k$ . As an example, we can

consider that  $C_1$ ,  $C_2$ , and  $C_3$  are respectively representing the  $R$ ,  $G$ , and  $B$  components in the RGB colour space. So  $P(I_t, C_k, i, j, )$  is representing the value of the colour component  $C_k$  for the pixel with coordinates  $i$  and  $j$  in frame  $I_t$ .

Some methods are dealing with histograms. So we define  $H(I_t, v)$  the number of pixels of the image  $I_t$  with an intensity equal to  $v$ , with  $v \in [0, V]$  where  $V$  is the maximum gray-level value. If we consider colour images, indexing methods can use several histograms, one for each colour component. We then use the notation  $H(I_t, C_k, v)$  to define the number of pixels with an intensity value of  $v$  for the colour component  $C_k$  in the image  $I_t$ .

Another common approach for video segmentation is to use block-sampled images. Let us note  $B$  the number of blocks  $b$  in each frame.

Finally, because a lot of methods are using some thresholds for shot change detection, we have also noted  $T$  some threshold fixed by the user. Several authors [11, 24] propose a learning procedure in order to use an appropriate threshold value.

As can easily be imagined from this introduction part, the works dealing with video sequence segmentation are quite numerous. We are reporting 94 entries in our bibliography. Of course some others exist but we think covering the main ways used to solve the problem. Even if the complexity of the methods is naturally increasing along time we have not chosen a chronological thread to present the various methods. Rather we have sorted them according to the basic elements they are relying on. We have organized them from the most simple, the pixel in the image to the most sophisticated ones, those that are using a combination of methods. More precisely we have distinguished 6 large categories characterized by the respective use of:

- pixel characterization,
- histogram of the frames,
- partition of the image in blocks,
- features,
- motion during the sequence,
- and combination of approaches.

## 2 Pixel-based methods

Shot change detection can be performed by comparing successive frames. The simplest way to compute the dissimilarity between two frames is to compare corresponding pixels from two successive images [61]. As we will see, some improvements of the initial pixel comparison have been proposed. First we present the methods considering two consecutive frames and then those that extend the study to a longer temporal interval.

## 2.1 Pixel comparison between two successive frames

One of the first method described in literature was from Nagasaka et al [61] in 1991. Shot changes are detected using a simple global interframe difference measure, defined as:

$$\text{Detection if: } \left( \left| \sum_{i=1}^X \sum_{j=1}^Y P(I_t, i, j) - \sum_{i=1}^X \sum_{j=1}^Y P(I_{t-1}, i, j) \right| \right) > T \quad (1)$$

resulting in  $O(\mathcal{P})$  operations per frame (as the second term of the difference has been already obtained after the processing of the previous frame  $I_{t-1}$ ).

Nagasaka et al [61] also introduced a shot change detection method based on pixel pair difference called template matching. For every two successive frames, differences of intensities are computed on pixels having the same spatial position in the two frames. Then the cumulated sum of differences is compared to a fixed threshold in order to determine if a shot change has been detected:

$$\text{Detection if: } \left( \sum_{i=1}^X \sum_{j=1}^Y |P(I_t, i, j) - P(I_{t-1}, i, j)| \right) > T \quad (2)$$

The number of operations per frame is equal to  $O(3\mathcal{P})$ . A colour version (of higher complexity  $O(9\mathcal{P})$ ) has also been presented:

$$\text{Detection if: } \left( \sum_{i=1}^X \sum_{j=1}^Y \sum_{k=1}^3 |P(I_t, C_k, i, j) - P(I_{t-1}, C_k, i, j)| \right) > T \quad (3)$$

A couple of years later, Zhang et al [91] were comparing the pixels of two successive frames on a boolean basis. The fact that pixels are different is noted:

$$D(I_t, I_{t-1}, i, j) = \begin{cases} 1 & \text{if } P(I_t, i, j) \neq P(I_{t-1}, i, j) \\ 0 & \text{otherwise} \end{cases} \quad (4)$$

for the gray-level case and requires one operation per couple of pixels. Definition is quite similar for colour images. In order to allow some variations on pixel intensities, a better (but more complex as it needs three operations instead of one) definition is given:

$$D(I_t, I_{t-1}, i, j) = \begin{cases} 1 & \text{if } |P(I_t, i, j) - P(I_{t-1}, i, j)| > T_D \\ 0 & \text{otherwise} \end{cases} \quad (5)$$

where  $T_D$  is considered as the tolerance value. The amount of different pixels is computed and is compared to a given threshold, which results in the detection or not of a shot change:

$$\text{Detection if: } \left( \sum_{i=1}^X \sum_{j=1}^Y D(I_t, I_{t-1}, i, j) \right) > T \quad (6)$$

resulting in complexity of  $O(2\mathcal{P})$  or  $O(4\mathcal{P})$  according to the condition used to compare pixels. In order to avoid false detections due to motion in the video sequence, they also

propose to smooth the images with a filter of size  $3 \times 3$  before computing the  $D$  values. The filter limits the effects due to noise and camera motion.

Several other statistical measures have been proposed in the literature [71]. The normalized difference energy and the normalized sum of absolute differences can be used for shot detection, as shown by the following equations:

$$\text{Detection if: } \left( \frac{\sum_{i=1}^X \sum_{j=1}^Y (P(I_t, i, j) - P(I_{t-1}, i, j))^2}{\left( \sum_{i=1}^X \sum_{j=1}^Y P(I_t, i, j)^2 \right) \left( \sum_{i=1}^X \sum_{j=1}^Y P(I_{t-1}, i, j)^2 \right)} \right) > T \quad (7)$$

$$\text{Detection if: } \left( \frac{\sum_{i=1}^X \sum_{j=1}^Y |P(I_t, i, j) - P(I_{t-1}, i, j)|}{\sum_{i=1}^X \sum_{j=1}^Y P(I_t, i, j) + \sum_{i=1}^X \sum_{j=1}^Y P(I_{t-1}, i, j)} \right) > T \quad (8)$$

These measures are respectively characterized by a complexity equal to  $O(5P)$  and  $O(4P)$ . Indeed, in both methods, the second part of the denominator has been obtained after processing the previous frame  $I_{t-1}$  and so does not need to be computed once again.

## 2.2 Pixel intensity time variation

The previous two frame study can be generalized by analysing variations of intensities through time. Taniguchi et al [79] label pixels with respect to the evolution of their intensities on several successive frames. The labels used are “constant”, “step( $I_t$ )”, “linear( $I_{t_1}, I_{t_2}$ )”, and “no label”. These labels represent respectively pixels with constant values, pixels with a change in value at frame  $I_t$ , pixels with a progressive change in value between frames  $I_{t_1}$  and  $I_{t_2}$ , and finally pixels with random values due to motion. Two boolean conditions  $\Theta_1(I_{t_1}, I_{t_2}, i, j)$  and  $\Theta_2(I_{t_1}, I_{t_2}, i, j)$  (needing respectively 4 and 6 operations per pixel) are introduced in order to define the constancy of a set of pixel values  $P(I_t, i, j)$  with  $t_1 \leq t \leq t_2$ :

$$\Theta_1(I_{t_1}, I_{t_2}, i, j) = \begin{cases} \text{true} & \text{if } \left( \max_{t_1 \leq t \leq t_2} P(I_t, i, j) - \min_{t_1 \leq t \leq t_2} P(I_t, i, j) \right) < T \\ \text{false} & \text{otherwise} \end{cases} \quad (9)$$

$$\Theta_2(I_{t_1}, I_{t_2}, i, j) = \begin{cases} \text{true} & \text{if } \left( \begin{array}{l} \max_{t_1 \leq t \leq t_2} (P(I_t, i, j) + (t - t_1)\theta_{t_1, t_2}) \\ - \min_{t_1 \leq t \leq t_2} (P(I_t, i, j) + (t - t_1)\theta_{t_1, t_2}) \end{array} \right) < T \\ \text{false} & \text{otherwise} \end{cases} \quad (10)$$

with  $\theta_{t_1, t_2}$  defined as:

$$\theta_{t_1, t_2} = \frac{|P(I_{t_1}, i, j) - P(I_{t_2}, i, j)|}{t_2 - t_1} \quad (11)$$



which requires 3 operations per pixel to be obtained since  $t_2 - t_1$  is a constant value, computed only once per couple of frames. These similarity conditions  $\Theta_1$  and  $\Theta_2$  are then used to determine the label  $L(I_{t_0}, I_{t_f}, i, j)$  of each pixel of a video sequence involving  $f + 1$  frames, using the following scheme:

$$L(I_{t_0}, I_{t_f}, i, j) = \begin{cases} \text{constant} & \text{if } \Theta_1(I_{t_0}, I_{t_f}, i, j) \\ \text{step}(I_t) & \text{if } \left( \begin{array}{l} \Theta_1(I_{t_0}, I_{t-1}, i, j) \\ \wedge \quad \Theta_1(I_t, I_{t_f}, i, j) \\ \wedge \quad \neg \Theta_1(I_{t-1}, I_t, i, j) \end{array} \right) \\ \text{linear}(I_{t_1}, I_{t_2}) & \text{if } \left( \begin{array}{l} \Theta_1(I_{t_0}, I_{t_1}, i, j) \\ \wedge \quad \Theta_1(I_{t_2}, I_{t_f}, i, j) \\ \wedge \quad \neg \Theta_1(I_{t_1}, I_{t_2}, i, j) \end{array} \right) \\ \text{no label} & \text{otherwise} \end{cases} \quad (12)$$

which can also be defined as:

$$L_{\text{label}}(I_{t_0}, I_{t_f}, i, j) = \begin{cases} 1 & \text{if } L(I_{t_0}, I_{t_f}, i, j) \text{ is of kind "label"} \\ 0 & \text{otherwise} \end{cases} \quad (13)$$

Quantities of pixels associated with each label are computed. Cuts (respectively dissolves) are detected thanks to the analysis of the ratio between quantity of pixels labelled “step” (respectively “linear”) and quantity of pixels labelled (*i.e.* with a label different of “no label”). A cut is detected at frame  $I_t$  if:

$$\text{Detection if: } \left( \frac{\sum_{i=1}^X \sum_{j=1}^Y L_{\text{step}(I_t)}(I_{t_0}, I_{t_f}, i, j)}{XY - \sum_{i=1}^X \sum_{j=1}^Y L_{\text{no label}}(I_{t_0}, I_{t_f}, i, j)} \right) > T \quad (14)$$

A dissolve is detected between frames  $I_{t_1}$  and  $I_{t_2}$  if:

$$\text{Detection if: } \left( \frac{\sum_{i=1}^X \sum_{j=1}^Y L_{\text{linear}(I_{t_1}, I_{t_2})}(I_{t_0}, I_{t_f}, i, j)}{XY - \sum_{i=1}^X \sum_{j=1}^Y L_{\text{no label}}(I_{t_0}, I_{t_f}, i, j)} \right) > T \quad (15)$$

Considering a number of operations per pixel respectively equal to at least 4 for equation (9), 6 for equation (10), 3 for equation (11), and 2 for equation (14) or (15), the overall complexity is then equal to  $O(15P)$ .

Lawrence et al [49] use evolution of temporal derivative of the pixel intensities as a criterion for shot change detection. First pixels with high spatial derivative are discarded in order to avoid motion effect. A pixel  $P(I_t, i, j)$  is considered if and only if the following condition holds:

$$\max(|P(I_t, i, j) - P(I_t, i - 1, j)|, |P(I_t, i, j) - P(I_t, i, j - 1)|) < T \quad (16)$$

This condition requires 6 operations per pixel to be verified. A convolution process involving remaining pixels and a gaussian mask is then performed to obtain temporal derivative of  $P(I_t, i, j)$ . Absolute values of these derivatives are summed up in order to define the distance measure (needing two operations per pixel) which will be analysed through time. Shot boundaries correspond to local maxima of this distance measure. False detections due to noise or motion are limited if the neighbourhood of the local maxima obtained previously are further analysed. Method complexity is equal to  $O(8\mathcal{P})$  without considering the gaussian filtering.

### 3 Histogram-based methods

The previous section was dedicated to pixel-based methods. It is also possible to compare two images based on global features instead of local features (pixels). Histogram is a global image feature widely used in image processing. The main advantage of histogram-based methods is their global aspect. So they are more robust to camera or object motion. The main drawback appears when we compare two different images having a similar histogram. It will often results in missing a shot change.

Different uses of the histogram can be distinguished. Some methods only compute differences between histograms and then the quality of the result is linked to the kind of histogram considered. A first extension is the use of weighted differences between histograms. Another approach consists in the definition of an intersection operator between histograms or the definition of different distances or similarity measures.

#### 3.1 Histogram difference

Tonomura et al [80] proposed a method based on gray-level histograms. Images are compared by computing a distance (of complexity  $O(3\mathcal{N})$ ) between their histograms, as shown in the following equation:

$$\text{Detection if: } \left( \sum_{v=0}^V |H(I_t, v) - H(I_{t-1}, v)| \right) > T \quad (17)$$

Nagasaka et al [61] propose a similar method using only 64 bins for colour histograms (2 bits for each colour component of the RGB space). Using the notation  $H_{64}(I_t, v)$ , the detection is defined by:

$$\text{Detection if: } \left( \sum_{v=0}^{63} |H_{64}(I_t, v) - H_{64}(I_{t-1}, v)| \right) > T \quad (18)$$

Gargi et al apply histogram difference to other colour spaces (HSV, YIQ, L\*a\*b\*, L\*u\*v\*, and Munsell) in [31]. More precisely, only non-intensity components are used (*i.e.* Hue and Saturation for HSV, I and Q for YIQ, a\* and b\* for L\*a\*b\*, u\* and v\* for L\*u\*v\*, and hue and chroma for the Munsell space). Shot change detection is then

defined by:

$$\text{Detection if: } \left( \sum_{k=1}^2 \sum_{v=0}^V |H(I_t, C_k, v) - H(I_{t-1}, C_k, v)| \right) > T \quad (19)$$

As it uses two colour components instead of only one, complexity is twice higher (*i.e.* equal to  $O(6N)$ ).

Pye et al [69] compute three histogram differences, considering separately the three colour components of the RGB space. The highest value is compared to a threshold for a shot change detection of complexity  $O(9N)$ :

$$\text{Detection if: } \left( \max_{k \in \{R, G, B\}} \sum_{v=0}^V |H(I_t, C_k, v) - H(I_{t-1}, C_k, v)| \right) > T \quad (20)$$

Ahmed et al [5] present several shot change detection algorithms using colour histograms. The first algorithm compares two frames using histograms computed on the Hue component  $C_H$ . So the detection needs  $O(4N)$  operations and can be represented by:

$$\text{Detection if: } \left( \frac{\sum_{v=0}^V |H(I_t, C_H, v) - H(I_{t-\Delta}, C_H, v)|}{\sum_{v=0}^V H(I_{t-\Delta}, C_H, v)} \right) > T \quad (21)$$

where  $\Delta$  is the temporal skip between two frames.

The second algorithm by Ahmed et al is based on reduced RGB space histograms. As in [61], histograms are composed of only 64 bins, using 2 bits for each colour component. The detection is done through a computation similar to the previously mentioned method:

$$\text{Detection if: } \left( \frac{\sum_{v=0}^{63} |H_{64}(I_t, v) - H_{64}(I_{t-\Delta}, v)|}{\sum_{v=0}^{63} H_{64}(I_{t-\Delta}, v)} \right) > T \quad (22)$$

resulting in a similar complexity.

O'Toole et al [66] detect shot changes using a cosine similarity measure computed between two histograms. First three 64 bin histograms representing respectively the Y, U, and V components are obtained from each frame. Next the three histograms are concatenated into a single one in order to get only one 192 bin histogram per frame. Then two successive frames are compared based on their histogram using a cosine similarity measure to perform shot change detection:

$$\text{Detection if: } \left( 1 - \frac{\sum_{v=0}^V (H_{YUV}(I_t, v) H_{YUV}(I_{t-1}, v))}{\sum_{v=0}^V H_{YUV}(I_t, v)^2 \sum_{v=0}^V H_{YUV}(I_{t-1}, v)^2} \right) > T \quad (23)$$

The method complexity is  $O(4\mathcal{N})$ . A similar work has been done by Cabedo et al in [18].

Chiu et al [21] rely their video segmentation on a genetic algorithm using colour histogram differences of complexity  $O(9\mathcal{N})$ . Possible shot boundaries are evaluated with similarity adjacency functions. In order to limit the optimization cost of these functions, a genetic algorithm is used instead of traditional methods. A video sequence is encoded as a string of binary values, 1 and 0 representing respectively the presence or not of a shot boundary in the current frame. The fitness function used in the algorithm is defined as a similarity adjacency function based on colour histogram differences. Finally crossover and mutation processes are derived from classical genetic algorithms in order to be adapted to video segmentation task.

Zhang et al [91] propose a method called twin comparison. Successive frames are compared using a histogram difference metric of complexity  $O(3\mathcal{N})$ . The difference values obtained are compared with two thresholds. Cuts are detected when difference is higher than a high threshold  $T_H$ . Possible starts of gradual transition are detected when difference is higher than a low threshold  $T_L$ . In this case, an accumulated difference is computed until the current difference is below  $T_L$ . Finally the accumulated difference is compared to the high threshold  $T_H$  for shot change detection. The two thresholds can be automatically set using standard deviation and mean of the interframe differences in the whole video sequence.

Li et al [54] use also a two step method, detecting successively the location of the end of the transition and its start. Frames are compared using the colour ratio histogram metric [47]. First two frames  $I_{t_1}$  and  $I_{t_2}$  (with  $t_2 = t_1 + \Delta$ ) are compared using this metric. While the difference is below a given threshold  $T$ ,  $t_2$  is set to  $t_2 + 1$ . When the difference is above  $T$ , the transition end has been obtained. In order to determine the transition start,  $t_1$  is set to  $t_2 - 1$ . The difference between frames  $I_{t_1}$  and  $I_{t_2}$  is then computed and compared to the threshold  $T$ . While the difference is below  $T$ ,  $t_1$  is set to  $t_1 - 1$ . When the difference is above  $T$ , the transition start has also been obtained.

Several other statistical measures have been reviewed in [71]. The quadratic histogram difference can be computed between histograms from two successive frames, whereas the Kolmogorov-Smirnov statistic is computed between cumulative histograms from two successive frames. These two measures are detailed below, using the notation  $H_C(I_t, v)$  to represent the cumulative histogram up to bin  $v$  for the frame  $I_t$ .

$$\text{Detection if: } \left( \sum_{v=0}^V \frac{(H(I_t, v) - H(I_{t-1}, v))^2}{(H(I_t, v) + H(I_{t-1}, v))^2} \right) > T \quad (24)$$

$$\text{Detection if: } \left( \max_{v \in [0, V]} (|H_C(I_t, v) - H_C(I_{t-1}, v)|) \right) > T \quad (25)$$

The two methods are characterized by a complexity respectively equal to  $O(6\mathcal{N})$  and  $O(3\mathcal{N})$ .

## 3.2 Weighted difference

In colour images, some colour components may have a bigger influence than others. So it is possible to detect shot changes by weighting the histograms of each colour

component depending on their importance [22]:

$$\text{Detection if: } \left( \sum_{k=1}^3 \sum_{v=0}^V \frac{L(I_t, C_k)}{L_{\text{mean}}(I_t)} |H(I_t, C_k, v) - H(I_{t-1}, C_k, v)| \right) > T \quad (26)$$

where  $L(I_t, C_k)$  and  $L_{\text{mean}}(I_t)$  are respectively the luminance for the  $k^{\text{th}}$  colour component of the frame  $I_t$  and the average luminance of the frame  $I_t$  considering all the colour components. The method complexity is equal to  $O(3\mathcal{P} + 9\mathcal{N})$ .

Zhao et al [92] use a learning procedure to determine the best weight values for weighted histogram difference computation. They first compute the original histogram difference with complexity  $O(3\mathcal{N})$  defined by equation (17). Then a learning step formulated as a minmax optimization problem is performed in order to select the best weights to use in weighted histogram differences. The detection process relies finally on the following equation which requires 12 operations per histogram bin:

$$\text{Detection if: } \left( \sum_{k=1}^3 \sum_{v=0}^V w(k, v) |H(I_t, C_k, v) - H(I_{t-1}, C_k, v)| \right) > T \quad (27)$$

where  $w(k, v)$  represents the best weight selected after the learning step. The overall complexity is then  $O(15\mathcal{N})$ .

Gargi et al presented in [31] a method based on the difference of average colours of a histogram, which can be as well considered as a histogram weighted difference. The shot change detection can then be represented by:

$$\text{Detection if: } \left( \sum_{k=1}^3 \left( \sum_{v=0}^V H(I_t, C_k, v)v - \sum_{v=0}^V H(I_{t-1}, C_k, v)v \right)^2 \right) > T \quad (28)$$

and requires  $O(6\mathcal{N})$  operations.

Another method by Gargi et al using colour histograms has also been described in [31]. More precisely, it uses a reference colour table as a frame difference measure. Reference colour table can be seen as a coarse quantization of RGB colour space into 27 different colour triples which are used as bins for a 3-D colour histogram  $H_{\text{ref}}$ . The shot change detection needs 5 operations per bin and can be represented by:

$$\text{Detection if: } \left( \sum_{v=0}^V w(v, t) \sqrt{(H_{\text{ref}}(I_t, v) - H_{\text{ref}}(I_{t-1}, v))^2} \right) > T \quad (29)$$

where the weight  $w(v, t)$  is defined as:

$$w(v, t) = \begin{cases} H_{\text{ref}}(I_{t-1}, v) & \text{if } (H_{\text{ref}}(I_t, v) > 0) \wedge (H_{\text{ref}}(I_{t-1}, v) > 0) \\ 1 & \text{otherwise} \end{cases} \quad (30)$$

and requires 3 operations per bin. The overall complexity is then equal to  $O(8\mathcal{N})$ .

### 3.3 Histogram intersection

Similarity between two images can also be evaluated thanks to histogram intersection. Histogram intersection is computed using different operators, for example a min function. In this case the computational cost is  $O(2\mathcal{N})$ . Similarity ratio belonging to interval  $[0, 1]$  is then compared to a given threshold. This comparison allows the detection of shot changes:

$$\text{Detection if: } \left( 1 - \frac{1}{XY} \sum_{v=0}^V \min(H(I_t, v), H(I_{t-1}, v)) \right) > T \quad (31)$$

where  $XY$  represents the number of pixels in frames processed.

Another version of the histogram intersection-based shot change detection method is defined in [31] using the following equation:

$$\text{Detection if: } \left( 1 - \frac{1}{XY} \sum_{v=0}^V \frac{\min(H(I_t, v), H(I_{t-1}, v))}{\max(H(I_t, v), H(I_{t-1}, v))} \right) > T \quad (32)$$

with a complexity equal to  $O(4\mathcal{N})$ .

Haering et al [36] apply histogram intersection defined in equation (31) to HSV (Hue, Saturation, Value) colour space, using 16 bins for Hue component and 4 bins each for Saturation and Value components.

An extension of [36] has been proposed by Javed et al in [43]. Hue is represented using only 8 bins. Instead of thresholding the histogram intersection of two successive frames, they compute the difference between two successive histogram intersection values and compare this derivative to a threshold.

### 3.4 Use of $\chi^2$ Test

Nagasaka et al have also proposed in [61] a 64 bin histogram comparison based on  $\chi^2$  test. The shot change detection is then defined by:

$$\text{Detection if: } \left( \sum_{v=0}^{63} \frac{(H_{64}(I_t, v) - H_{64}(I_{t-1}, v))^2}{H_{64}(I_t, v)} \right) > T \quad (33)$$

with the assumption  $H_{64}(I_t, v) \neq 0$ . If this assumption does not hold, we use the following equation instead:

$$\text{Detection if: } \left( \sum_{v=0}^{63} \frac{(H_{64}(I_t, v) - H_{64}(I_{t-1}, v))^2}{H_{64}(I_{t-1}, v)} \right) > T \quad (34)$$

with the assumptions  $H_{64}(I_{t-1}, v) \neq 0$  and  $H_{64}(I_t, v) = 0$ . This method is considered as more efficient than simple histogram comparison-based methods. Its complexity is equal to  $O(5\mathcal{N})$ .

A modification has been proposed by Dailianas et al in [22] where the detection is represented by:

$$\text{Detection if: } \left( \sum_{v=0}^V \frac{(H(I_t, v) - H(I_{t-1}, v))^2}{\max(H(I_t, v), H(I_{t-1}, v))} \right) > T \quad (35)$$

with a similar complexity.

Gunsel et al [35] perform a K-means clustering algorithm to determine the location of shot boundaries. Successive frames are compared using  $\chi^2$  test or histogram difference on YUV histograms, resulting in a complexity equal to  $O(15N)$ . Every interframe difference value is classified into shot change or non-shot change.

### 3.5 Similarity measures between normalized histograms

Several measures computed on normalized histograms have been reviewed by Ren et al in [71] and by Kim and Park in [46]. Using the notation  $H_N(I_t, v)$  to represent the probability of intensity  $v$  in the frame  $I_t$ , cross entropy, divergence, Kullback Liebler distance, and Bhattacharya distance are respectively defined as:

$$\text{Detection if: } \left( \sum_{v=0}^V \left( H_N(I_t, v) \log \frac{H_N(I_t, v)}{H_N(I_{t-1}, v)} \right) \right) > T \quad (36)$$

$$\text{Detection if: } \left( \begin{array}{l} \sum_{v=0}^V \left( H_N(I_t, v) \log \frac{H_N(I_t, v)}{H_N(I_{t-1}, v)} \right) \\ + \sum_{v=0}^V \left( H_N(I_{t-1}, v) \log \frac{H_N(I_{t-1}, v)}{H_N(I_t, v)} \right) \end{array} \right) > T \quad (37)$$

$$\text{Detection if: } \left( \begin{array}{l} \sum_{v=0}^V \left( H_N(I_t, v) \log \frac{H_N(I_t, v)}{H_N(I_{t-1}, v)} \right) \\ + \sum_{v=0}^V \left( (1 - H_N(I_t, v)) \log \frac{1 - H_N(I_t, v)}{1 - H_N(I_{t-1}, v)} \right) \end{array} \right) > T \quad (38)$$

$$\text{Detection if: } \left( -\log \left( \sum_{v=0}^V \sqrt{H_N(I_t, v) H_N(I_{t-1}, v)} \right) \right) > T \quad (39)$$

with complexities respectively equal to  $O(4N)$ ,  $O(8N)$ ,  $O(11N)$ , and  $O(3N)$ .

All these methods are based on a uniform process all over the image. The heterogeneity present within a frame led to use block-based methods.

## 4 Block-based methods

Block sampling of the video frames can be performed in order to increase the quality of shot change detection but also to decrease the computation time. Once block representation has been obtained from original images, it is possible to perform some algorithms derived from pixel or histogram-based methods presented previously. Use of

blocks allows a processing which is intermediate, between local level like pixel-based methods and global level as histogram-based methods. Main advantage of block-based methods is their relative insensitivity to noise and camera or object motion. We have distinguished between several approaches all working on blocks.

#### 4.1 Block similarity

Kasturi et al [45] perform a similarity test on block-sampled images. Like in pixel-based methods, pairs of blocks (with same spatial coordinates) from two successive frames are compared. The similarity is based on block features like mean and variance, which can be computed on a complete frame considering a complexity of  $O(\mathcal{P} + \mathcal{B})$  and  $O(2\mathcal{P} + 3\mathcal{B})$ . The likelihood rate  $L$  (of complexity  $O(9\mathcal{B})$ ) is defined for a block  $b$  as:

$$L(I_t, I_{t-1}, b) = \frac{\left( \frac{\sigma_{t,b}^2 + \sigma_{t-1,b}^2}{2} + \left( \frac{\mu_{t,b} - \mu_{t-1,b}}{2} \right)^2 \right)^2}{\sigma_{t,b}^2 \sigma_{t-1,b}^2} \quad (40)$$

where  $\mu_{t,b}$  and  $\sigma_{t,b}^2$  are respectively the mean and the variance of block  $b$  pixel values in image  $I_t$ . Then thresholded values  $L_D$  of  $L$  are defined by the equation:

$$L_D(I_t, I_{t-1}, b) = \begin{cases} 1 & \text{if } L(I_t, I_{t-1}, b) > T_D \\ 0 & \text{in other cases} \end{cases} \quad (41)$$

where  $T_D$  is considered as a tolerance value. A detection is obtained when:

$$\text{Detection if: } \sum_{b=1}^B c_b L_D(I_t, I_{t-1}, b) > T \quad (42)$$

where  $c_b$  is used to give more or less importance to block  $b$ . Most of the time  $c_b$  is set to 1 for all the blocks. Overall complexity is estimated to  $O(3\mathcal{P} + 15\mathcal{B})$  considering required operations for estimation of the block mean, variance, likelihood rate, and thresholding, and the final cost of the detection. This likelihood ratio can also be used directly on full-frames, as proposed in [71].

Another well-known measure involving variance is the Yakimovsky likelihood ratio which can be applied also on blocks or frames directly [71]. For each block this ratio is computed as:

$$L'(I_t, I_{t-1}, b) = \frac{\left( \sigma_{\{t,t-1\},b}^2 \right)^2}{\sigma_{t-1,b}^2 \sigma_{t,b}^2} \quad (43)$$

where  $\sigma_{t,b}^2$  and  $\sigma_{t-1,b}^2$  represent the variances of the pixel intensity values in the frames  $I_t$  and  $I_{t-1}$  considering a block  $b$ . The notation  $\sigma_{\{t,t-1\},b}^2$  is used to denote the variance of the pooled data from both frames for a block  $b$ . Knowing  $\sigma_{t,b}^2$  and  $\sigma_{t-1,b}^2$ , computation of  $\sigma_{\{t,t-1\},b}^2$  needs 4 operations per block. When  $\mu_{t,b}$  is not available, computation of  $\sigma_{t,b}^2$  is characterized by a cost of  $O(3\mathcal{P} + 4\mathcal{B})$ . The overall complexity  $O(3\mathcal{P} + 13\mathcal{B})$  is obtained by adding the cost of equation (43) (*i.e.* 3 operations per block).



Freund statistic can also be used to detect shot changes. Distance measure is then defined by:

$$L''(I_t, I_{t-1}, b) = \frac{\mu_{t,b} - \mu_{t-1,b}}{\sqrt{\frac{\sigma_{t,b}^2 + \sigma_{t-1,b}^2}{XY}}} \quad (44)$$

resulting in a complexity of  $O(3\mathcal{P} + 11\mathcal{B})$  obtained by replacing cost of equation (40) by cost of Freund statistic, *i.e.*  $O(5\mathcal{B})$ .

Lee et al [52] perform shot change detection using block differences computed in the HSV colour space. First RGB images are converted to HSV in order to avoid camera flashes. Then the mean values of Hue and Saturation components are computed for each block (with a cost of  $O(\mathcal{P} + \mathcal{B})$  for each colour component). Two successive blocks are compared using these mean values:

$$D_H(I_{t_1}, I_{t_2}, b) = |\mu(I_{t_1}, b, C_H) - \mu(I_{t_2}, b, C_H)| \quad (45)$$

$$D_S(I_{t_1}, I_{t_2}, b) = |\mu(I_{t_1}, b, C_S) - \mu(I_{t_2}, b, C_S)| \quad (46)$$

where  $\mu(I_t, b, C_k)$  is the mean of a bloc  $b$  in the frame  $I_t$  considering the colour component  $C_k$ .  $D_H$  and  $D_S$  represent respectively differences for Hue and Saturation colour component. These two distances need 2 operations per block each and are used to determine if each block has changed:

$$D(I_{t_1}, I_{t_2}, b) = \begin{cases} 1 & \text{if } (D_H(I_{t_1}, I_{t_2}, b) > T_H) \vee (D_S(I_{t_1}, I_{t_2}, b) > T_S) \\ 0 & \text{otherwise} \end{cases} \quad (47)$$

which requires 3 operations per block. Finally the ratio between the number of changed blocks and the total number of blocks is compared to another threshold in order to detect shot changes:

$$\text{Detection if: } \frac{1}{B} \sum_{b=1}^B D(I_t, I_{t-\Delta}, b) > T \quad (48)$$

where  $\Delta$  represents the temporal skip used in the shot change detection process. The overall complexity of this method is  $O(2\mathcal{P} + 10\mathcal{B})$ .

## 4.2 Histogram comparison

Swanberg et al present in [78] a method detecting shot changes thanks to the comparison of colour histograms computed on the blocks of the images noted  $H(I_t, b, C_k, v)$ . The detection process is then defined as:

$$\text{Detection if: } \left( \sum_{k=1}^3 \sum_{b=1}^B \sum_{v=0}^V \frac{(H(I_t, b, C_k, v) - H(I_{t-1}, b, C_k, v))^2}{H(I_t, b, C_k, v) + H(I_{t-1}, b, C_k, v)} \right) > T \quad (49)$$

and has a computational cost of  $O(15\mathcal{N}\mathcal{B} + 3\mathcal{B})$ .

Nagasaka et al [61] extend their histogram comparison to images divided in 4x4 blocks. Every pair of blocks from two successive frames is compared using the  $\chi^2$  test

on 64 bin histograms:

$$\chi(I_t, b) = \left( \sum_{v=0}^{63} \frac{(H_{64}(I_t, b, v) - H_{64}(I_{t-1}, b, v))^2}{H_{64}(I_t, b, v)} \right) \quad (50)$$

which requires for each bloc 4 operations per histogram bin. The values obtained are then sorted in an ascending way and the 8 lowest are kept. The sum of these values is computed and compared to a threshold to detect shot changes:

$$\text{Detection if: } \left( \sum_{b=1}^8 \chi_s(I_t, b) \right) > T \quad (51)$$

where the  $\chi_s$  values represent ascending sorted values of  $\chi$  (*i.e.* for  $b \in [1, 16]$  we have  $\chi_s(I_t, 1) \leq \chi_s(I_t, b) \leq \chi_s(I_t, 16)$ ). Considering equations (50) and (51), the overall complexity of the method is  $O(4NB + 0.5B)$ .

Ueda et al [82] proposed to use the rate of correlation change instead of the magnitude of correlation change proposed in [61]. Each value  $\chi(I_t, b)$  obtained from a pair of blocks is compared to a threshold. The detection depends on the number of significant values  $\chi(I_t, b)$  instead of the sum of the highest  $\chi(I_t, b)$ , resulting in a complexity equal to  $O(4NB + 2B)$ .

Ahmed et al [5] propose also a block-based version of their method using the 6 most significant RGB bits as described in equation (22). They compare histograms computed on blocks instead of global histograms. The sum of the histogram differences obtained for each block is computed and compared to a predefined threshold in order to detect shot changes, as shown in:

$$\text{Detection if: } \left( \sum_{b=1}^B \frac{\sum_{v=0}^{63} |H_{64}(I_t, b, v) - H_{64}(I_{t-\Delta}, b, v)|}{\sum_{v=0}^{63} H_{64}(I_{t-\Delta}, b, v)} \right) > T \quad (52)$$

where  $\Delta$  represents the temporal skip between two successive frames to be analysed. The computation cost is also equal to  $O(4NB + 2B)$ .

Ahmed et al proposed in [4] an improved version of their algorithm described previously. Instead of comparing two frames considering a fixed temporal skip, the method is based on an adaptive temporal skip. First, two images  $I_{t_1}$  and  $I_{t_2}$  are compared according to equation (52). Then if the difference is greater than a threshold,  $t_2$  is replaced by  $\frac{t_1+t_2}{2}$  and the frames are again compared. If the difference is still greater than the threshold,  $t_1$  is also set to  $\frac{t_1+t_2}{2}$  (considering the current values of  $t_1$  and  $t_2$ ) and frames are compared. This process is repeated until  $t_1 + 1 = t_2$  which represents a shot change between frames  $t_1$  and  $t_2$ .

Lee et al introduce in [50] a selective HSV histogram comparison algorithm. First, pixels are classified with respect to their colour level. If a pixel is characterized by high values for V and S, it is classified into a discrete colour using H component. Otherwise the classification is based on the intensity (or gray-level) value. For a given

pixel  $P(I_t, b, i, j)$ , two complementary states are defined:

$$S_{hue}(I_t, b, i, j) = \begin{cases} 1 & \text{if } \left( \begin{array}{l} (P(I_t, b, C_S, i, j) > T_S) \\ \wedge (P(I_t, b, C_V, i, j) > T_V) \end{array} \right) \\ 0 & \text{otherwise} \end{cases} \quad (53)$$

$$S_{gray}(I_t, b, i, j) = 1 - S_{hue}(I_t, b, i, j) \quad (54)$$

Computation of these states requires 4 operations per pixel. For each block two selective histograms  $H_{hue}(I_t, b)$  and  $H_{gray}(I_t, b)$  are then computed in a classical way considering the two states previously defined. The notation  $H_{hue}(I_t, b, v)$  (respectively  $H_{gray}(I_t, b, v)$ ) represents the number of pixels in block  $b$  of frame  $I_t$  with state  $S_{hue}$  (resp.  $S_{gray}$ ) equal to 1 and with hue (resp. insensity or gray-level) value equal to  $v$ . These histograms are used for shot change detection:

$$\text{Detection if: } \left( \sum_{b=1}^B \left( \begin{array}{l} \sum_{v=0}^V |H_{hue}(I_t, b, v) - H_{hue}(I_{t-1}, b, v)| \\ + \sum_{v=0}^V |H_{gray}(I_t, b, v) - H_{gray}(I_{t-1}, b, v)| \end{array} \right) \right) > T \quad (55)$$

resulting in a complexity equal to  $O(6\mathcal{NB} + 2\mathcal{B} + 4\mathcal{P})$ .

Bertini et al [13] compare in the HSI colour space histograms of successive frames divided in 9 blocks. In order to improve robustness to change in lighting conditions, the Intensity component is not used. The detection can then be represented by the following equations:

$$D_{HS}(I_t, I_{t+1}, b) = \sum_{k \in \{H, S\}} \sum_{v=0}^V (H(I_t, b, C_k, v) - H(I_{t+1}, b, C_k, v)) \quad (56)$$

$$D'_{HS}(I_t) = \sum_{b=1}^B D_{HS}(I_t, I_{t+1}, b) - \sum_{b=1}^B D_{HS}(I_{t-1}, I_t, b) \quad (57)$$

$$\text{Detection if: } \left( \begin{array}{l} ((D'_{HS}(I_t) > 0) \wedge (D'_{HS}(I_{t+1}) < 0)) \\ \vee ((D'_{HS}(I_t) < 0) \wedge (D'_{HS}(I_{t+1}) > 0)) \end{array} \right) \quad (58)$$

The two distances  $D_{HS}(I_t, I_{t+1}, b)$  and  $D'_{HS}(I_t)$  require respectively 2 operations per bin per block per colour and one operation per block. The overall complexity is then  $O(4\mathcal{NB} + \mathcal{B})$ . In order to improve the detection, a minimum temporal distance is defined between two successive cuts.

Chahir et al [19] based their method on histogram intersection computed on frames divided in 24 blocks. The colour space used in their method is  $L^*u^*v^*$ . For each block, color histogram intersection is computed between two successive frames requiring 12 operations per bin. A comparison with a threshold allows to determine whether the block has been changed or not. The number of changed blocks is then compared to another threshold in order to detect a shot change, resulting in an overall complexity equal to  $O(12\mathcal{NB} + 2\mathcal{B})$ .

### 4.3 Combination of histogram differences and likelihood rate

This method proposed by Dugad in [25] is based on two successive steps to detect cuts and other transitions. Shot changes are detected using successively histogram difference and likelihood ratio. In this method three thresholds have to be set. Histogram difference step (whose complexity is  $O(3\mathcal{N})$ ) is defined as in equation (17) and is compared with two thresholds. The difference is first compared to a high threshold in order to avoid false alarms. If it is lower than this threshold, it is then compared to a low threshold. If it is higher than this low threshold, the final decision is taken by computing likelihood ratio values. In this case, the two frames to be compared are divided in 64 blocks and the 16 central blocks are kept. For each block, the likelihood ratio is computed between the block  $P(I_t, b)$  and the blocks  $P(I_{t-1}, b')$  where  $b'$  belongs to the neighbourhood of  $b$ , and the minimum of the likelihood value is kept. Then a mean of the 16 minimum likelihood ratios is computed and is compared with the third threshold, which may result in a shot change detection. The respective costs for computation of the mean and variance of each block and estimation of the likelihood ratio in a  $3 \times 3$  neighbourhood are  $O(3\mathcal{P} + 4\mathcal{B})$  and  $O(72\mathcal{B})$ . Computation of the minimum values for each block and the mean value on all block requires respectively 7 and 1 operations per block. As we are using only a quarter of the blocs, we divide the number of operations per block by 4. The overall complexity is then  $O(3\mathcal{P} + 3\mathcal{N} + 21\mathcal{B})$ .

### 4.4 Use of neighbourhood colour ratio

Adjeroh et al [2] compare two successive frames using neighbourhood colour ratios. A local averaging step (with a cost of  $O(5\mathcal{P} + \mathcal{B})$ ) is first performed in order to obtain one value  $P'(I_t, b)$  per block:

$$P'(I_t, b) = \prod_{i=2}^{X-1} \prod_{j=2}^{Y-1} \frac{1}{4P(I_t, b, i, j)} \left( \begin{array}{l} P(I_t, b, i-1, j) + P(I_t, b, i+1, j) \\ + P(I_t, b, i, j-1) + P(I_t, b, i, j+1) \end{array} \right) \quad (59)$$

Pairs of blocks from two different frames are then compared using this measure:

$$D'(I_t, I_{t-\Delta}, b) = 1 - \min \left( \frac{P'(I_t, b)}{P'(I_{t-\Delta}, b)}, \frac{P'(I_{t-\Delta}, b)}{P'(I_t, b)} \right) \quad (60)$$

which requires 4 operations per block and where  $\Delta$  represents the temporal skip. Shot changes are finally detected if the number of significant  $D'$  values for all selected blocks is higher than a fixed threshold, or if the average value of  $D'$  is higher than another threshold. These two conditions need respectively 2 and 1 operations per block to be verified. The overall complexity of this method is then  $O(5\mathcal{N}\mathcal{B} + 8\mathcal{B})$ .

### 4.5 Evolution of block dissimilarity

Shot changes can also be detected by analysing the evolution of block dissimilarity. Demarty et al [23] compute locally a distance criterion (of cost  $O(9\mathcal{P})$ ) in RGB colour space between blocks of two successive images. Result obtained consists in distance

values between the two images for every block. Then the sum of these values is computed (which requires 3 operations per block) and an evolution curve of this sum is built. This evolution curve is filtered using a top-hat morphological operation and is finally compared with a threshold in order to detect shot changes. The complexity is equal to  $O(9\mathcal{P} + 3\mathcal{B})$ .

Lefèvre et al [53] proposed a method using HSV colour space on block-sampled images in order to avoid false detection due to illumination effects. A value is obtained for each block from Hue and Saturation components with a cost in  $O(2\mathcal{P} + 5\mathcal{B})$ . Then a block-based difference (requiring 3 operations per block) is computed between two frames based on the block values. This difference is tracked through time, as well as its derivative. Analysis of this derivative allows cut detection, whereas the initial (non-derivated) difference values are used to initialize a cumulative sum computation of the derived values. This allows detection of gradual transitions. This method is characterized by a computational cost of  $O(2\mathcal{P} + 8\mathcal{B})$ .

#### 4.6 Temporal and spatial subsampling

Xiong et al [85] propose to subsample the video sequence in both space and time. An abrupt change is detected between two frames  $I_t$  and  $I_{t+\Delta}$  if:

$$\text{Detection if: } \left( \sum_{b \in B'} D_\mu(I_t, I_{t+\Delta}, b) \right) > T \quad (61)$$

where  $B'$  represents a set of a priori selected blocks and  $D_\mu(I_{t_1}, I_{t_2}, b)$  is defined as:

$$D_\mu(I_{t_1}, I_{t_2}, b) = \begin{cases} 1 & \text{if } |\mu_{t_1, b} - \mu_{t_2, b}| > T_\mu \\ 0 & \text{in other cases} \end{cases} \quad (62)$$

The two equations need respectively 1 and 3 operations per block. As the block mean computation is linked with a cost of  $O(\mathcal{P} + \mathcal{B})$ , the overall complexity is equal to  $O(\mathcal{P} + 5\mathcal{B})$ . Gradual transitions are detected using an edge-based frame-to-frame difference measure. If a shot change is detected, a binary search is performed reducing  $\Delta$  to determine the exact shot boundaries. The method proposed is called ‘‘Net Comparison’’ and has also been tested in HSV colour space.

### 5 Feature-based methods

All the methods we have already presented were using features, but they can be qualified of trivial features. Here we are considering more sophisticated ones. We consider:

- the moments computed on the image,
- the contour lines extracted from the image,
- some feature points extracted using Hough Transform,
- the planar points,

- colour transition,
- modelling of the video transition effects,
- the use of some decision process as bayesian methods,
- features computed from classical statistical approaches,
- and the use of Hidden Markov Models.

## 5.1 Moment invariants

Arman et al [12] use moment invariants combined with histogram intersection to detect shot changes. Moment invariants have properties such as invariance to scale change, rotation, and translation. The moments of a frame  $I_t$  are defined as:

$$m_{p,q} = \sum_{i=1}^X \sum_{j=1}^Y i^p j^q P(I_t, i, j) \quad (63)$$

In [22], shot changes are detected thanks to the computation of the usual Euclidean distance between two frames using a vector composed of the first three moment invariants, defined as:

$$\vec{\Phi} = \begin{pmatrix} m_{2,0} + m_{0,2} \\ (m_{2,0} - m_{0,2})^2 + 4m_{1,1}^2 \\ (m_{3,0} - 3m_{1,2})^2 + (3m_{2,1} - m_{0,3})^2 \end{pmatrix} \quad (64)$$

Considering equation (63), all moments used in equation (64) require 3 operations per pixel, excepting  $m_{1,2}$  and  $m_{2,1}$  which need 4 operations per pixel. The detection can be finally defined as:

$$\text{Detection if: } \left( \left\| \vec{\Phi}(I_t) - \vec{\Phi}(I_{t-1}) \right\|^2 \right) > T \quad (65)$$

resulting in an overall complexity of  $O(23\mathcal{P})$ .

## 5.2 Edges

Zabih et al in [90] use edge extraction to detect shot changes. Global motion compensation is performed on successive frames. Next, edges are extracted using Canny algorithm and dilated. Normalized proportions of entering edges and exiting edges are then computed (with a cost of  $O(3\mathcal{P})$  each) using the following equations:

$$P(C_{out}, I_t) = 1 - \frac{\sum_{i=1}^X \sum_{j=1}^Y E(I_{t-1}, i + \alpha_{t-1,t}, j + \beta_{t-1,t}) E_d(I_t, i, j)}{\sum_{i=1}^X \sum_{j=1}^Y E(I_{t-1}, i, j)} \quad (66)$$

$$P(C_{in}, I_t) = 1 - \frac{\sum_{i=1}^X \sum_{j=1}^Y E_d(I_{t-1}, i + \alpha_{t-1,t}, j + \beta_{t-1,t}) E(I_t, i, j)}{\sum_{i=1}^X \sum_{j=1}^Y E(I_{t-1}, i + \alpha_{t-1,t}, j + \beta_{t-1,t})} \quad (67)$$

where  $E$  and  $E_d$  are respectively the contour image and its dilated version, and  $(\alpha_{t-1,t}, \beta_{t-1,t})$  represents the global motion translation vector between the two successive images  $I_t$  and  $I_{t-1}$ . Then a dissimilarity measure  $ECF(I_t)$  called edge change fraction is computed by:

$$ECF(I_t) = \max(P(C_{out}, I_t), P(C_{in}, I_t)) \quad (68)$$

Finally this value is compared to a threshold to detect shot changes:

$$\text{Detection if: } ECF(I_t) > T \quad (69)$$

resulting in a complexity at least equal to  $O(26\mathcal{P})$  when considering edge detection requires 20 operations per pixel.

Smeaton et al proposed an evolution of the previous method in [77] where the detection is based on the evolution of the edge change fraction on several frames instead of the analysis of this dissimilarity measure on only one frame. Detection can then be defined by:

$$\text{Detection if: } (ECF(I_t) - ECF(I_{t-1})) > T \quad (70)$$

Lienhart [56] also uses edge information to perform dissolve detection. First edges extracted with the Canny edge detector are confronted with two thresholds to determine weak and strong edges:

$$E_w(I_t, i, j) = \begin{cases} E(I_t, i, j) & \text{if } T_w \leq E(I_t, i, j) \leq T_s \\ 0 & \text{in other cases} \end{cases} \quad (71)$$

$$E_s(I_t, i, j) = \begin{cases} E(I_t, i, j) & \text{if } T_s \leq E(I_t, i, j) \\ 0 & \text{in other cases} \end{cases} \quad (72)$$

where  $T_w$  and  $T_s$  are respectively the lowest and highest thresholds for detecting weak and strong edges.  $E_w$  and  $E_s$  represent the weak and strong edge images and need respectively 2 and 1 operations per pixel to be computed. Then the edge-based contrast  $EC$  is obtained for a frame  $I_t$  according to the equation:

$$EC(I_t) = 1 + \frac{\sum_{i=1}^X \sum_{j=1}^Y E_s(I_t, i, j) - \sum_{i=1}^X \sum_{j=1}^Y E_w(I_t, i, j) - 1}{\sum_{i=1}^X \sum_{j=1}^Y E_s(I_t, i, j) + \sum_{i=1}^X \sum_{j=1}^Y E_w(I_t, i, j) + 1} \quad (73)$$

with a cost of  $O(2\mathcal{P})$ . Finally dissolves are detected when the current value of  $EC$  is a local minimum. The overall complexity  $O(25\mathcal{P})$  is obtained by adding edge detection cost.

Yu et al [87] use edge information to detect gradual transitions. Cuts are first detected using a histogram difference measure (of cost  $O(3\mathcal{P})$ ) computed between two

successive subsampled frames. Then a second pass is necessary for detecting gradual transitions. For every frame  $I_t$  between two successive cuts at time  $t_1$  and  $t_2$ , the number  $Q_E(I_t)$  of edge points present in the image is computed (considering 21 operations per pixel) and temporally smoothed. Then for every local minima  $Q_E(I_{t'})$  which is below a predefined threshold, a search of the two closest local maxima  $Q_E(I_{t'_1})$  and  $Q_E(I_{t'_2})$  is performed with  $t'_1 < t' < t'_2$ . A fade out effect is detected between  $t'_1$  and  $t'$  if:

$$\text{Detection if: } \left( \frac{1}{t' - t'_1} \sum_{t=t'_1}^{t'} |Q_E(I_t) - Q_E(I_{t-1})| \right) \in [T_{\text{low}}^{\text{fade out}}, T_{\text{high}}^{\text{fade out}}] \quad (74)$$

Similarly, a fade in effect is detected between  $t'$  and  $t'_2$  if:

$$\text{Detection if: } \left( \frac{1}{t'_2 - t'} \sum_{t=t'}^{t'_2} |Q_E(I_t) - Q_E(I_{t-1})| \right) \in [T_{\text{low}}^{\text{fade in}}, T_{\text{high}}^{\text{fade in}}] \quad (75)$$

For dissolve effect detection, a new measure called double chromatic difference is computed for every frame belonging to the interval  $[t'_1, t'_2]$ :

$$DCD(I_t) = \sum_{i=1}^X \sum_{j=1}^Y \left| \frac{1}{2} P(I_{t'_1}, i, j) + \frac{1}{2} P(I_{t'_2}, i, j) - P(I_t, i, j) \right| \quad (76)$$

which requires 6 operations per pixel. If the frame number  $t_{\min}$  corresponding to the global minimum value  $DCD(I_{t_{\min}})$  is close enough to  $t'$  ( $|t' - t_{\min}| < \epsilon$  with  $\epsilon$  defined as a small number), a dissolve effect is assumed to be found between frames  $I_{t'_1}$  and  $I_{t'_2}$ . The overall complexity of this method is equal to  $O(3N + 27P)$ .

Heng and Ngan [41] also propose a method based on edge information. They introduce the notion of edge object, considering the pixels close to the edge. Occurrences of every edge object are matched on two successive frames. Shot changes are detected using ratio between the amount of edge objects persistent over time and the total amount of edge objects.

Nam and Tewfik [62] propose a coarse-to-fine shot change detection method based on wavelet transforms. Image sequences are first temporally subsampled. Frames processed are also spatially reduced using a spatial 2-D wavelet transform. Intensity evolution of pixels belonging to coarse frames is analysed using a temporal 1-D wavelet transform. Sharp edges define possible shot change locations. Video frames around these locations are further processed at full-rate. Temporal 1-D wavelet transform is applied again on the full-rate video sequence. Edge detection is also performed on every coarse frame and the number of edge points is computed on a block-based basis. Difference between two successive frames is computed using the number of edge points for each block. True shot boundaries are located on sharp edges in the 1-D wavelet transform and high values of interframe difference considering block-based amount of edge points. An extension to wipe transitions detection has been proposed in [63].



### 5.3 Feature points

Ardebilian et al [11] detect shot changes by comparing between feature points extracted from two successive images. They use Hough transform to extract feature points. Success or not of the feature points matching between two successive frames results directly in cut detection.

### 5.4 Planar points

Silva et al [76] perform shot change detection using spatio-temporal representation of an image sequence. A video sequence  $\mathcal{V}$  is represented in  $\mathbb{R}^4$  as an hypersurface:

$$\mathcal{V} = \{i, j, t, P(I_t, i, j)\} \quad (77)$$

The amount of planar points in every frame is considered as the measure for detecting cuts. We recall that planar points are points contained in a flat neighbourhood of the hypersurface. For a given frame  $I_t$ , planar points are determined using the characteristic polynomial coefficients of the Hessian matrix of  $P(I_t, i, j)$ . Then the percentage of planar points is computed. A cut is detected (in a four frame interval) when this value is greater than three times the temporal variance of the percentage of planar points. The overall complexity  $O(51\mathcal{P})$  is obtained by adding the cost associated with the Hessian matrix ( $O(27\mathcal{P})$ ), the three characteristic polynomial coefficients ( $O(18\mathcal{P})$ ), the classification of a point as planar ( $O(5\mathcal{P})$ ), and the final decision ( $O(\mathcal{P})$ ).

### 5.5 Colour transitions

Sanchez et al [74] compare between two successive frames using colour histograms computed on specific regions. These regions are defined from the most significant colour transitions of the image, considered as high values of its multispectral gradient and computed with Sobel approximation. Colour histograms are compared between regions of two successive frames to determine the coherence of the region through time. Shot changes are finally detected if the amount of changed regions is above a given threshold.

### 5.6 Transition modelling

Some shot changes are created from production effects. These transitions can be modelled explicitly with mathematical tools in order to be detected. Several methods using these assumptions are presented below.

Hampapur et al [37] model several kinds of fades and wipes with mathematical functions. Knowing the two last shots and their respective durations, it is possible to estimate the duration of the current shot. Detection of shot changes can rely on a constancy measure (of cost  $O(\mathcal{P})$ ) defined for frame  $I_t$  as:

$$C(I_t) = \frac{\sum_{i=1}^X \sum_{j=1}^Y S_{\text{binary}}(I_t, i, j)}{\sigma_t(1 + |X_c(I_t) - \frac{X}{2}| + |Y_c(I_t) - \frac{Y}{2}|)} \quad (78)$$

where  $\sigma_t$  represents the standard deviation of pixel intensities in frame  $I_t$  computed with a cost of  $O(3\mathcal{P})$ . The binary state  $S_{\text{binary}}(I_t, i, j)$  of a pixel is defined as:

$$S_{\text{binary}}(I_t, i, j) = \begin{cases} 1 & \text{if } P(I_t, i, j) \neq 0 \\ 0 & \text{otherwise} \end{cases} \quad (79)$$

requiring 1 operation per pixel and the components of the centroid of image  $I_t$  are noted  $X_c(I_t)$  and  $Y_c(I_t)$ . They are computed as:

$$X_c(I_t) = \frac{\sum_{i=1}^X \sum_{j=1}^Y iP(I_t, i, j)}{\sum_{i=1}^X \sum_{j=1}^Y P(I_t, i, j)} \quad (80)$$

$$Y_c(I_t) = \frac{\sum_{i=1}^X \sum_{j=1}^Y jP(I_t, i, j)}{\sum_{i=1}^X \sum_{j=1}^Y P(I_t, i, j)} \quad (81)$$

and need together 5 operations per pixel. The overall complexity is then  $O(10\mathcal{P})$ .

Adami et al [1] perform dissolve detection applying a model of dissolve effects on frame histograms. For every frame, two specific histograms are computed:

$$\bar{H}(I_t) = H_{\frac{1}{2}}(I_t) * H_{\frac{1}{2}}(I_t) \quad (82)$$

$$H_{\Delta}(I_t) = H_{\frac{1}{2}}(I_{t-\Delta}) * H_{\frac{1}{2}}(I_{t+\Delta}) \quad (83)$$

where  $H_{\frac{1}{2}}(I_t)$  represents the histogram of the frame  $I_t$  scaled by half,  $\Delta$  is a fixed parameter, and the operator  $*$  figures a convolution. Both operations are characterized by a cost of  $O(\mathcal{N}^2)$ . It is then possible to compute a dissimilarity measure using these histograms:

$$R(I_t) = \frac{\chi^2(H(I_t), \bar{H}(I_t))}{\chi^2(H(I_t), H_{\Delta}(I_t))} - 1 \quad (84)$$

where the  $\chi^2$  operator (of complexity  $O(5\mathcal{N})$ ) is computed between two histograms. Maxima values of  $R(I_t)$  indicate dissolve locations:

$$\text{Detection if: } R(I_t) > R(I_{t'}) \quad \forall t' \text{ in the neighbourhood of } t \quad (85)$$

The overall complexity is then  $O(2\mathcal{N}^2 + 10\mathcal{N})$ .

Aigrain et al [6] detect shot changes using the assumption of the absence of important motion in the different shots. Their method is based on a motion differential model and uses a density function estimation as the difference between two images. First, two successive images are reduced spatially and normalized using histogram equalization. Then the histogram of pixel-pair difference is computed and is simplified to two values, which are respectively the amount of differences belonging to the interval  $[128, 255]$  computed on normalized images and the amount of differences belonging to the interval  $[1, 40]$  computed on non-normalized images. Both values require 4 operations per

pixel in order to be computed. Cut and gradual transition detections are respectively based on local maxima of the first and second value described previously. The method complexity is equal to  $O(8P)$ .

Lienhart [56] relies on fade modelling from [37] to perform fade detection. The proposed algorithm uses the standard deviation of pixel intensities as an estimation of the scaling function introduced in fade effects. First, all monochrome frames are located as they are fade start or end candidates. These frames are characterized by a standard deviation  $\sigma(I_t)$  close to zero. Fades are then detected by searching in both temporal directions for a linear increase in the standard deviation.

Alattar [10] bases also the shot change detection on variance of pixel intensities. Fades are detected using a two steps scheme. First local minimum values of the second order difference of the pixel intensity spatial variance time series are obtained. Then a test is performed to determine whether the first order difference of the pixel intensity mean is relatively constant in the neighbourhood of the local variance minimum or not. In the positive case, a fade is assumed to be found. A similar method [9] has been proposed for dissolve.

Truong et al [81] combine approaches from [56] and [10]. First, all monochrome frames are detected. Then only monochrome frames which are next to a local minimum of the intensity variance second order difference are kept. Finally a test is performed on the first order difference of the mean intensity. If this value is constant through time and other conditions are satisfied, a fade has been detected. The other conditions correspond to comparison between thresholds and the absolute value of the first order difference and the intensity variance of the first or last frame. Dissolve detection is performed using the evolution of the variance first order difference through time. This difference value should be monotonically increasing from a negative value up to a positive value. So zero crossings are used to locate dissolve frames.

Fernando et al [30] use also mean and variance of the luminance signal to determine fade and dissolve transitions locations. For every frame, the mean and the variance of the luminance is computed. The ratio between second temporal derivative of the variance to the first temporal derivative of the mean is then compared between two successive frames. Shot changes are located when this ratio is constant through time.

All the approaches based on effect modelling using variance and mean of successive frames are characterized by a cost of  $O(3P)$ .

## 5.7 Bayesian approaches

Vasconcelos et al propose in [83] a segmentation method using a Bayesian model of the editing process. For each frame a local activity measure is computed based on a tangent distance. In order to detect shot changes, this measure is compared (following a Bayesian framework) to an adaptive threshold, depending on the *a priori* duration of a shot and on the time elapsed between the previous shot change and the current frame  $I_t$ .

Hanjalic et al [40] use also a statistical framework for the shot change detection, which is modelled as a probability minimization problem of the average detection error. Detection criteria are linked with visual content discontinuity (based on motion compensation) and knowledge about shot length distribution.

Han et al [39] base their detection on gray-level or colour histogram differences computed between successive frames using equation (17) or equation (19). A filtering step combining an average clipping operation and a subsequent local convolution is used to improve the shot change detection. The evolution curve of the filtered histogram difference value is analysed and decision for the detection of a shot change is taken following a Bayesian framework. Detections of a cut or a gradual effect are respectively linked with the presence in the evolution curve of a rectangular or triangular shape.

## 5.8 Statistical approaches

Yilmaz et al [86] use Principal Coordinate System to perform shot change detection on RGB frames. First, image rows are concatenated in order to obtain only one row vector per colour component for each frame. We use the notations  $V(I_t, C_R)$ ,  $V(I_t, C_G)$  and  $V(I_t, C_B)$  for the row vectors associated with the Red, Green and Blue components. Then the 3x3 covariance matrix  $M(I_t)$  of the RGB colour space is computed following:

$$M(I_t) = \begin{pmatrix} V(I_t, C_R) \\ V(I_t, C_G) \\ V(I_t, C_B) \end{pmatrix} \begin{pmatrix} V(I_t, C_R)^T & V(I_t, C_G)^T & V(I_t, C_B)^T \end{pmatrix} \quad (86)$$

which requires 18 operations per pixel. Next the vector representing the principal axis is selected and noted  $V_{\lambda_{\max}}(I_t)$ . We recall this vector is the eigenvector associated with the maximum eigenvalue  $\lambda_{\max}$  of the covariance matrix. Finally two successive frames are compared with respect to the angle between their respective principal axes following the equation:

$$\text{Detection if: } \left( 1 - \frac{V_{\lambda_{\max}}(I_t)^T V_{\lambda_{\max}}(I_{t-1})}{\|V_{\lambda_{\max}}(I_t)\| \|V_{\lambda_{\max}}(I_{t-1})\|} \right) > T \quad (87)$$

so the actual complexity is equal to  $O(18\mathcal{P})$ .

Han et al [38] use also a principal component analysis to perform shot change detection. First frames are subsampled and represented as column vectors. Then successive frames are grouped in a temporal window. The mean vector  $\mu$  and the empirical covariance matrix  $M$  of this window are computed. Then the unique set of orthonormal eigenvectors of  $M$  and their associated eigenvalues are obtained. Each frame in the window is then projected onto the  $K$  eigenvectors corresponding to the  $K$  largest eigenvalues. Finally shot changes are detected using the temporal variations of angle and length of the  $K$  first principal components.

Li et al [55] based their algorithm on the computation of joint probability images between frames. They use a spatio-temporal representation of the successive joint probability images obtained in order to detect shot changes. First a joint probability image is computed (with a cost of  $O(2\mathcal{P})$ ) between two frames, which consists in the frequency of the co-occurrences of intensity or chrominance values. Two similar images  $I_{t_1}$  and  $I_{t_2}$  will be characterized by a joint probability image  $J(I_{t_1}, I_{t_2})$  composed of non-zero values on the diagonal. A distance measure is then defined between two frames using

the joint probability image:

$$D_J(I_{t_1}, I_{t_2}) = 1 - \frac{\sum_{(i,j) \in \mathcal{J}'} J(I_{t_1}, I_{t_2}, i, j)}{\sum_{(i,j) \in \mathcal{J}} J(I_{t_1}, I_{t_2}, i, j)} \quad (88)$$

where  $\mathcal{J}$  and  $\mathcal{J}'$  represent respectively the set of all pixels  $(i, j)$  in the joint probability image and the set of all pixels near the diagonal line with a given tolerance  $\delta$  (*i.e.*  $\mathcal{J}' = \{(i, j) : |i - j| < \delta\}$ ). As  $J$  and  $J'$  contain about  $N^2$  and  $N$  values, the computation of  $D_J$  has a cost of  $O(N^2 + N)$ . If the value  $D_J$  obtained is higher than a fixed threshold, several algorithms are used in order to confirm the presence of a shot change. Dissolve effects are detected using histogram intersection performed on spatio-temporal representations of joint probability images. This method is characterized by a complexity equal to  $O(2P + N^2 + N)$ .

Gong and Liu [34] perform shot change detection using the Singular Value Decomposition. Every frame is divided in  $3 \times 3$  blocks on which a 3-D colour histogram composed of 125 bins is computed. A vector of 1125 components is then obtained for every frame. The video sequence is represented by a matrix which is processed by a singular value decomposition algorithm. The  $K$  largest singular values are kept and are used to compute a similarity measure between two frames. Detection of a shot change is done by comparing the similarity computed between the two frames  $I_{t_1}$  and  $I_{t_2}$  with a low and a high threshold. If the similarity measure is below the low threshold, no shot change has been detected. On the contrary, if the measure is higher than the high threshold, a shot change is assumed to be found. In the last case (*i.e.* the similarity measure is between the two thresholds), a refinement step is performed involving frames between  $I_{t_1}$  and  $I_{t_2}$ .

## 5.9 Hidden Markov models

Eickeler et al [26] use Hidden Markov Models to perform video indexing. Some of the classes represent shot boundary frames. Several features are defined to describe each frame, but only some of them characterize shot boundary frames:

$$d_1(I_t) = \min \left( \frac{\sum_{i=1}^X \sum_{j=1}^Y |P(I_t, i, j) - P(I_{t-1}, i, j)|}{XY}, \frac{\sum_{i=1}^X \sum_{j=1}^Y |P(I_{t+1}, i, j) - P(I_{t-2}, i, j)|}{XY} \right) \quad (89)$$

$$d_2(I_t) = \sum_{v=0}^V \left( -\text{median} \left( \begin{array}{c} |H(I_t, v) - H(I_{t-1}, v)| \\ |H(I_{t-1}, v) - H(I_{t-2}, v)|, \\ |H(I_t, v) - H(I_{t-1}, v)|, \\ |H(I_{t+1}, v) - H(I_t, v)| \end{array} \right) \right) \quad (90)$$

$$d_3(I_t) = \sum_{i=1}^X \sum_{j=1}^Y \frac{P(I_t, i, j) - P(I_{t-1}, i, j)}{P(I_t, i, j) - 0.5P(I_{t-1}, i, j) + P(I_{t+1}, i, j)} \quad (91)$$

where  $d_1(I_t)$ ,  $d_2(I_t)$ , and  $d_3(I_t)$  represent respectively the intensity of motion, the median filtered intensity of difference histograms, and the ratio between the difference pixel and the difference from interpolated pixel. These three measures need respectively 6 operations per pixel, 6 operation per histogram bin, and 5 operations per pixel, resulting in an overall complexity equal to  $O(11\mathcal{P} + 6\mathcal{N})$ . After a learning step using the Baum-Welch algorithm, segmentation is performed using the Viterbi algorithm.

A similar approach using Hidden Markov Models has been proposed by Boreczky in [15]. The model is based on image, audio, and motion features. Classification of a frame into a shot boundary class is done using only the image feature of the frame. This feature is defined as a luminance 64 bin histogram difference similar to the one described in equation (17).

## 6 Motion-based methods

As the nature of motion is usually continuous in a video sequence, it can also be used as a criterion to detect shot changes. Based on this fact, several approaches using motion information were proposed in the literature. We review here methods based on global (or camera) motion, motion vectors, optical flow, and correlation in the frequency domain.

### 6.1 Global motion

Cherfaoui and Bertin [20] detect shot changes in two steps. First the global motion parameters are estimated using an affine transformation model. The estimated motion is then used to classify a shot as fixed, pan, or zoom. If the motion is not coherent through time, a shot change is assumed to be found.

Bouthemy et al [16] based their detection on a dominant multiresolution motion estimation. This estimation uses a global 2-D parametric model composed of 6 parameters. Once the dominant motion has been estimated, a coefficient  $\omega_{i,j}$  is also available for every pixel  $(i, j)$ . It represents the coherence of the pixel with the dominant motion estimated. Using a predefined threshold, it is possible to define the set of dominant motion-coherent pixels in each frame. The evolution of the set size through time is analysed in order to detect shot changes.

Zugaj et al [93] extend the previous method to wipe detection. Here only pixels which are non-coherent with the estimated dominant motion are considered. For each frame, two histograms are computed based on the number of non-coherent pixels along horizontal and vertical axes. For every couple of frames, absolute differences between corresponding histograms are computed and result in two other histograms. The correlation between two successive absolute difference histograms is then measured along the two axes. If one of the two correlation values exceeds a predefined threshold, an horizontal or vertical wipe is detected.

Mann et al [60] proposed a method where global motion estimation is performed using a 8 parameter model. They define “video orbits” as collection of pictures starting from one picture and applying all possible projective coordinate transformations to that picture using the 8 motion parameters. Two frames belonging to the same scene will lie

in the same orbit or nearly so. So shot changes are detected when the distance between the orbits of successive frames is higher than a threshold.

Complexity of global or dominant motion computation is usually estimated at  $O(20\mathcal{P})$  when considering only 2-D translations.

## 6.2 Motion vectors

Akutsu et al [8] use a motion smoothness measure to detect shot changes. The indexing method uses subsampled video sequences and processes only one frame every  $k$  frames. Then the selected frame is divided into  $8 \times 8$  blocks and each block is matched to a block in the next chosen frame. Motion vector is estimated thanks to the closest matched neighbouring block, which is also used to compute the value of the correlation coefficient. An interframe similarity measure can be defined as the average of these correlations. Another measure called motion smoothness is defined as the ratio between the number of blocks which have significantly moved and the displacement of these blocks. Shot changes are finally detected in presence of local extrema in the correlation and motion smoothness ratio values.

Shahraray [75] proposed a similar method. Subsampled video sequences are used and every frame is divided in 12 blocks. A research is performed to match each block from one frame to the most similar block (in a spatial neighbourhood) in the next frame. Motion vector and correlation value are computed in a way similar to the Akutsu et al method. Main difference with the previous method is the use of a nonlinear digital order statistic filter. Correlation values are sorted into ascending order and the first values and their respective motion vectors are used for the computation of an average value which is considered as a similarity measure. As in [8], a local temporal extremum in the similarity measure means shot change detection. A motion-controlled temporal filter is used to avoid false detection due to motion.

Liu et al [58] based their method on motion compensated images obtained from motion vector information. First motion vectors of frame  $I_{t-1}$  are used to create a motion compensated version  $I'_t$  of the frame  $I_t$ . The next step is luminance normalization. The motion compensated frame  $I'_t$  is normalized in order to get the same energy as the original frame  $I_t$ . Normalized image is noted  $I''_t$ . The original frame  $I_t$  is then compared to the two modified frames  $I'_t$  and  $I''_t$  using a modified version of the  $\chi^2$  test applied on their histograms. The result  $\chi(I_t, I'_t)$  is compared to an adaptive threshold in order to detect cuts. Fade detection is based on the comparison between  $\chi(I_t, I'_t)$  and  $\chi(I_t, I''_t)$  which are the two histogram differences computed previously.

When motion vectors are obtained using a fast block-matching technique (as the three step search method or the 2-D log search method) and considering a search size of  $16 \times 16$  pixels, the complexity is estimated at  $O(75\mathcal{P})$ .

## 6.3 Optical Flow

Fatemi et al [28] use optical flow as information to detect shot changes. First the video sequence is divided into overlapping subsequences, defined as 3 consecutive frames and a fourth predicted frame. Every frame is then divided into  $B$  blocks, which are predicted from the first frame to the second one, and from the second frame to the third

one. Finally a set of 3 matching blocks from the first three frames is used for block prediction into the last frame. If the block prediction does not correctly estimate the content of the last frame, a shot change is assumed to be found.

Optical flow computation is usually characterized in the literature by a cost of  $O(105\mathcal{P})$ .

#### 6.4 Frequency domain correlation

Porter et al [68] propose a technique inspired by motion-based algorithms. Correlation between two successive frames is computed and used as a shot change detection measure. In order to compute the interframe correlation, a block-based approach working in the frequency domain is taken. Frames are divided into blocks of  $32 \times 32$  pixels. Every block in a frame  $I_{t-1}$  is matched with a neighbouring block in frame  $I_t$  by first computing the normalized correlation between blocks and then seek and located the correlation coefficient with the largest magnitude. The normalized correlation is computed in the frequency domain instead of the spatial domain to limit computation time. The average correlation is then obtained for a couple of frames. Shot changes are detected in presence of local minima of this value.

Complexity of motion estimation obtained by computing the correlation in the frequency domain is assumed to be similar to complexity of fast block-matching algorithms, *i.e.*  $O(75\mathcal{P})$ .

### 7 Combination of several methods

Even if most of shot change detection methods working in the uncompressed domain can be categorized under pixel-, histogram-, block-, feature-, and motion-based methods, some cannot. In this section we will present some methods which do not enter in the previously defined categories.

Two systems have been proposed by Quénot et al in [70]. The first one is dedicated to cut detection and based on colour histogram comparison by  $\chi^2$  test [61] followed by a simplified version of the contour analysis algorithm [90] where the motion compensation step is replaced by a dilation of the contour points. The two steps require respectively 5 operations per histogram bin and 26 operations per pixel, so the overall cost is equal to  $O(26\mathcal{P} + 5\mathcal{N})$ . In the second method, motion compensation is performed using optical flow. A direct analysis of the resulting images is done in order to detect cuts, whereas first and second temporal derivatives of the images are compared for progressive transition detection.

Gauch et al [33] detect shot changes thanks to the combination of three image features: the average brightness of each video frame, the change in pixel values and the change in colour distribution from two different frames. Shot changes are detected if at least one of the two following conditions holds:

$$\sum_{i=1}^X \sum_{j=1}^Y P(I_t, i, j) < T_1 \quad (92)$$



$$\left( \begin{array}{l} \left( \sum_{i=1}^X \sum_{j=1}^Y |P(I_t, i, j) - P(I_{t-\Delta}, i, j)| > T_2 \right) \\ \wedge \left( \sum_{v=0}^{255} |H_{256}(I_t, v) - H_{256}(I_{t-\Delta}, v)| > T_3 \right) \end{array} \right) \quad (93)$$

where  $\Delta$  represents the temporal skip and  $H_{256}$  the histograms quantified to 256 uniformly distributed colours. The two conditions are respectively characterized by a cost of  $O(\mathcal{P})$  and  $O(3\mathcal{P} + 3\mathcal{N})$ , resulting in an overall complexity of  $O(4\mathcal{P} + 3\mathcal{N})$ . Thresholds  $T_1$ ,  $T_2$ , and  $T_3$  are initialized manually at the beginning of the video using *a priori* knowledge and updated dynamically during the video analysis using statistical information.

Yusoff et al [89] perform shot change detection by combining 5 algorithms. The methods used are the average intensity measurement [37], the euclidean distance [84], the histogram comparison [91], the likelihood ratio [45], and the motion estimation method characterized respectively by a cost of  $O(3\mathcal{P})$ ,  $O(2\mathcal{B})$ ,  $O(3\mathcal{P} + 9\mathcal{N})$ ,  $O(3\mathcal{P} + 15\mathcal{B})$ , and  $O(75\mathcal{P})$ . The overall complexity is then  $O(85\mathcal{P} + 9\mathcal{N} + 17\mathcal{B})$ . The last algorithm uses prediction error computed as the sum of absolute differences between the original frame and the reconstructed frame built using motion information. Final decision is obtained using a majority vote (there should be at least 3 algorithms detecting a shot change to validate this hypothesis).

Sabata et al [73] perform fusion of multiple cues to parse video. The features used are from different origins. One is linked to colour distribution, another is a set of responses to several texture filters, some tracking information is also involved, and finally edges in spatio-temporal volumes are computed. Colour feature is analysed through the use of weighted histogram, defined for a temporal window  $[t_1, t_2]$  with  $t_1 < t_2$  as:

$$H(I_{t_1}, I_{t_2}, v) = \frac{1}{V} \sum_{t=t_1}^{t_2} w(t) H(I_t, v) \quad (94)$$

with the weight  $w(t)$  computed using the one dimensional Gaussian function  $G(t)$  multiplied by a normalization factor. Colour feature for the frame  $I_t$  is then defined as:

$$D^{colour}(I_t) = \sum_{v=0}^V \frac{(H(I_{t-\Delta}, I_t, v) - H(I_t, I_{t+\Delta}, v))^2}{H(I_{t-\Delta}, I_t, v) + H(I_t, I_{t+\Delta}, v)} \quad (95)$$

Texture feature is analysed in a similar way using a set of 12 Gabor filters instead of colour histogram. The 12 Gabor filters are computed using 4 values for  $\theta$  and 3 for  $\lambda$ , which represent respectively the orientation and the wavelength of the filter. The next feature is linked to results of feature tracking, where a score is assigned to the tracking in every frame. It is computed by weighting the contribution of each feature from the last frame to the current frame using the weight  $1 - e^{-\frac{\delta_i}{k}}$  where the constant  $k$  determines the sensitivity to the history of the track and  $\delta_i$  represents the number of frames in the past through which the  $i^{th}$  feature was tracked. The distance measure between two sets of frames  $[I_{t-\Delta}, I_t]$  and  $[I_t, I_{t+\Delta}]$  is then computed using difference between their average track scores and their missed tracks ratios. The last feature used is linked to spatio-temporal volumes. These volumes are built thanks to the projection

of the video data along the  $(x, t)$  and the  $(y, t)$  planes, as in [67]. For a given sequence, edges perpendicular to the  $t$  axis may indicate shot changes. So an average value of the fraction of pixels covered by the horizontal edges for every set of frames to be analysed is computed and is used as a probability measure of the existence of a shot change. Once all the features are obtained, the final decision is taken according to the result of a fusion process performed using a Bayesian network.

Naphade et al [64] detect shot changes using a K-means clustering algorithm using two metrics per frame couple. The two metrics are respectively based on histogram difference and pixel intensity difference defined in sections 2.1 and 3.1. Then all couples of successive frames are classified, based on these two metrics and using a K-means clustering algorithm with two clusters. The cluster characterized by higher values represents frames containing shot changes. A final step is necessary to eliminate false positive. Frames with a local maximum value for the histogram difference metric are selected as shot change boundaries. The complexity has been evaluated to  $O(4\mathcal{P} + 9\mathcal{N})$ .

Oh et al [65] propose a three step method based on a Gaussian pyramid representation of the background area of images. Background area is defined as a concatenation of three regions located respectively on the left, on the top, and on the right of the image. The first step, called pixel matching, is based on the computation of the pixel difference ratio, defined as:

$$\left( \frac{\sum_{i=1}^X \sum_{j=1}^Y D(I_t, I_{t-1}, i, j)}{XY} \right) \geq 0.10 \quad (96)$$

where  $D(I_t, I_{t-1}, i, j)$  is defined in equation (4). If this condition is satisfied, there may be a shot change. So in a next step a matching is performed between signs of two successive frames. A sign is defined as the pixel on the top of the gaussian pyramid. The matching can be represented by:

$$\left( \frac{1}{256} \max_{k \in \{R, G, B\}} \left( \begin{array}{c} \text{sign}^{BA}(I_t, C_k) \\ - \text{sign}^{BA}(I_{t-1}, C_k) \end{array} \right) \right) \geq 0.10 \quad (97)$$

where  $\text{sign}^{BA}(I_t, C_k)$  denotes the background area sign value of the  $C_k$  colour component for the frame  $I_t$ . If this second condition is satisfied, the detection process continues with a final step called background tracking. This is done by comparing signatures of two successive frames. Signatures are 1-D arrays extracted from background areas using Gaussian pyramid. The comparison is done by shifting the two arrays in opposite directions. In each matching step, only the overlapping pixels in the matching window are compared to determine the maximum continuous match, which is the maximum amount of consecutive pixels matched in the two frames. This measure is finally compared to a certain threshold in order to determine the presence of a shot change.

Ferman et al [29] extend the K-means based-method proposed in [35]. The unsupervised clustering process uses, for each couple of frames, two features which are filtered histogram difference (derivated from equation (17)) and filtered pair-wise pixel

comparison (derivated from equations (5) and (6)). For a given distance measure  $D(I_t)$ , the filtering is done in two steps:

$$D'(I_t) = \max \left( D(I_t) - \underset{t' \in [t-\Delta_1, t+\Delta_1]}{\text{mean}} D(I_{t'}) \right) \quad (98)$$

$$D''(I_t) = \max \left( D'(I_t) - \underset{t' \in [t-\Delta_2, t+\Delta_2]}{\text{median}} D'(I_{t'}) \right) \quad (99)$$

where  $\Delta_1$  and  $\Delta_2$  are temporal window lengths which have to be set. This method has a complexity similar to [64], *i.e.*  $O(4\mathcal{P} + 9\mathcal{N})$ .

Lee et al [51] also use a K-means clustering algorithm to perform shot change detection. Every couple of frames is represented by a vector of two components which are linked to pixel-based and histogram-based distance measures. The measures used are normalized versions of those described in equations (2) and (17). A vector is then composed of two features:

$$\begin{pmatrix} \frac{1}{255XY} \sum_{i=1}^X \sum_{j=1}^Y |P(I_t, i, j) - P(I_{t-1}, i, j)| \\ \frac{1}{2XY} \sum_{v=0}^V |H(I_t, v) - H(I_{t-1}, v)| \end{pmatrix} \quad (100)$$

which need respectively 3 operations per pixel and 3 operations per histogram bin. The complexity is then equal to  $O(3\mathcal{P} + 3\mathcal{N})$ . Every couple of frames  $(I_{t-1}, I_t)$  is classified using a 2-class K-means clustering algorithm applied on pairs of frames from  $(I_{t_0}, I_{t_0+1})$  to  $(I_{t_f-1}, I_{t_f})$  where  $I_{t_0}$  and  $I_{t_f}$  represent respectively the first and final frames of the video sequence. The two classes represent the set of couples of frames where a cut is present and its complement.

## 8 Complexity results

We reviewed in previous sections a great number of methods for segmentation of un-compressed video sequences. In order to compare the different methods, we have to select a criterion. Quality evaluation, as detailed in section 1.3, depends directly on the selection of optimal parameters for each method. So we decide to compare reviewed methods based on a complexity point of view rather than a quality one. As we are focusing on real time video parsing, this criterion is also of highest importance.

In this section, we will present complexity of reviewed methods. More precisely, theoretical and experimental results will be given.

### 8.1 Theoretical complexity

We present in tables 1 to 6 results of complexity computation for the different categories of video segmentation methods presented in this paper. Complexity computation results are given for most of the methods reviewed here. However complexity was

not computed for some methods providing insufficient available information, as for example video segmentation methods dealing with wavelet transform, Hough transform, or motion compensation.

Methods with complexity annotated by \* are characterized by higher complexity because of the explicit use of colour. Other methods were considered in a gray-level framework. Indeed, colour-based methods dealing with a subsampled colour space (for example 64 bin colour histogram-based methods) are considered as gray-level for complexity computation. For several methods, the complexity computed does not include some specific processing. So the actual complexities of the concerned methods are in fact higher than values given here. We use the notation + in this case.

We also have to notice that complexity of methods dealing with histograms (mainly those in table 2) does not include the cost of histogram computation. It is then necessary to add a cost of  $O(\mathcal{P})$  to every method using histograms.

Block-based methods (table 3) may not be characterized by a lower complexity. It is mainly due to extraction of block features which need a processing of all pixels in the image. Some of these methods (like [53]) can be used to segment compressed video sequences using for example the DC term as a block average value. In this case, the real complexity is much lower than using uncompressed frames. The same remark applies to several motion-based methods which detect shot changes using motion vectors. These vectors can be directly available in compressed video sequences whereas they have to be computed in uncompressed video sequences.

The theoretical estimations given here have been verified by actual experiments.

## 8.2 Experimental complexity

We estimated theoretically the complexity of most of the reviewed methods. In order to verify these estimations, an important part of reviewed methods were implemented. Experimental results are also given in tables 1 to 6.

As the computation time depends on the programming language and the computing architecture used, results are given in a relative form rather an absolute one. Simple interframe difference [61] being the fastest method, it is used as a reference computation time.

Tests have been performed using gray-level or RGB colour video sequences. Images are composed of  $192 \times 144$  pixels. The implementation of the reviewed methods has been performed using Matlab programming environment. As we compare methods on a relative basis, the absolute computation time is not critical. So we use Matlab rather than C or C++ in order to limit implementation time.

Due to the Matlab implementation, some apparently strange results have to be explained. All methods performing a colour space conversion (chrominance histogram difference, hue histogram difference, cosine similarity measure, block difference in the HSV colour space, selective HSV histogram comparison, and HSV block dissimilarity) are characterized by higher computation time as this kind of operation is quite slow using Matlab. Block-based methods were implemented using Matlab `blkproc` and `colfilt` functions to avoid loops, which perform really quite slow in Matlab. As `colfilt` performs faster than `blkproc` (ratios observed experimentally are between 135 % and 175 %), it was used as often as possible. However, even if we use

adequate Matlab functions, computation time of block-based methods stays high. Finally, motion-based methods are characterized by very variable computation times, depending on the use or not of Matlab optimized functions. As motion vectors are obtained using the block-matching technique which requires several nested loops, methods [8, 75, 58] using motion vectors show actual computation time worse than theoretical complexity.

## 9 Conclusion

We presented in this paper a great number of methods for segmentation of uncompressed video sequences. We focused on the complexity of these methods, contrary to other previously published reviews. This criterion is particularly important for two main reasons. First, uncompressed video contains a huge quantity of data, which is not the case of compressed video. So video parsing method dealing with uncompressed video will be characterized by their computational intensive aspect. Secondly, when dealing with real-time video segmentation, we have to use methods known for their low computational cost in order to process video at frame rate.

Results of this study show that the best method for real-time segmentation of uncompressed video sequences should be selected considering efficiency and computational cost. Simple methods like interframe difference is one of the fastest method but it is characterized by a poor quality. On the opposite, feature- or motion-based methods or methods combining several algorithms are more robust and lead to better quality, but they are also known to be greedy in computational resources. Another important criterion to choose a method is the way of computing threshold values and the number of thresholds. Future work will deal with the study of this criterion.

## References

- [1] N. Adami and R. Leonardi. Identification of editing effect in image sequences by statistical modelling. In *Symposium on Picture Coding*, pages 157–160, Portland, OR, April 1999.
- [2] D.A. Adjeroh, M.C. Lee, and C.U. Orji. Techniques for fast partitioning of compressed and uncompressed video. *International Journal of Multimedia Tools and Applications*, 4(2):225–243, March 1997.
- [3] G. Ahanger and T.D.C. Little. A survey of technologies for parsing and indexing digital video. *Journal of Visual Communication and Image Representation*, 7(1):28–43, March 1996.
- [4] M. Ahmed and A. Karmouch. Video segmentation using an opportunistic approach. In *Multimedia Modeling*, pages 389–405, Ottawa, Canada, October 1999.
- [5] M. Ahmed, A. Karmouch, and S. Abu-Hakima. Key frame extraction and indexing for multimedia databases. In *Vision Interface Conference*, pages 506–511, Trois-Rivières, Canada, May 1999.

- [6] P. Aigrain and P. Joly. The automatic real-time analysis of film editing and transition effects and its applications. *Computer and Graphics*, 18(1):93–103, January/February 1994.
- [7] P. Aigrain, H.J. Zhang, and D. Petkovic. Content-based representation and retrieval of visual media: A state-of-the-art review. *International Journal of Multimedia Tools and Applications*, 3(3):179–202, November 1996.
- [8] A. Akutsu, Y. Tonomura, H. Hashimoto, and Y. Ohba. Video indexing using motion vectors. In *SPIE Conference on Visual Communications and Image Processing*, volume 1818, pages 1522–1530, Boston, MA, November 1992.
- [9] A.M. Alattar. Detecting and compressing dissolve regions in video sequences with dvi multimedia image compression algorithm. In *IEEE International Symposium on Circuits and Systems*, volume 1, pages 13–16, Chicago, IL, May 1993.
- [10] A.M. Alattar. Detecting fade regions in uncompressed video sequences. In *International Conference on Acoustics, Speech, and Signal Processing*, pages 3025–3028, Munich, Germany, April 1997.
- [11] M. Ardebilian, X. Tu, and L. Chen. Robust 3d clue-based video segmentation for video indexing. *Journal of Visual Communication and Image Representation*, 11(1):58–79, March 2000.
- [12] F. Arman, R. Depommier, A. Hsu, and M.Y. Chiu. Content-based browsing of video sequences. In *ACM International Conference on Multimedia*, pages 97–103, San Francisco, CA, November 1994.
- [13] M. Bertini, A. Del Bimbo, and P. Pala. Content based indexing and retrieval of tv news. *Pattern Recognition Letters*, 22(5):503–516, April 2001.
- [14] J.S. Boreczky and L.A. Rowe. Comparison of video shot boundary detection techniques. *Journal of Electronic Imaging*, 5(2):122–128, April 1996.
- [15] J.S. Boreczky and L.D. Wilcox. A hidden markov model framework for video segmentation using audio and image features. In *International Conference on Acoustics, Speech, and Signal Processing*, volume 6, pages 3741–3744, Seattle, WA, May 1998.
- [16] P. Bouthemy, M. Gelgon, and F. Ganansia. A unified approach to shot change detection and camera motion characterization. *IEEE Transactions on Circuits and Systems for Video Technology*, 9(7):1030–1044, October 1999.
- [17] R. Brunelli, O. Mich, and C.M. Modena. A survey on the automatic indexing of video data. *Journal of Visual Communication and Image Representation*, 10(2):78–112, June 1999.
- [18] X.U. Cabedo and S.K. Bhattacharjee. Shot detection tools in digital video. In *Workshop on Non-linear Model Based Image Analysis*, pages 121–126, Glasgow, Scotland, July 1998.

- [19] Y. Chahir and L. Chen. Automatic video segmentation and indexing. In *SPIE Conference on Intelligent Robots and Computer Vision*, volume 3837, pages 345–356, Boston, MA, September 1999.
- [20] M. Cherfaoui and C. Bertin. Temporal segmentation of videos: a new approach. In *SPIE Conference on Digital Video Compression: Algorithms and Technologies*, volume 2419, pages 38–47, San Jose, CA, February 1995.
- [21] P. Chiu, A. Girgensohn, W. Polak, E. Rieffel, and L. Wilcox. A genetic algorithm for video segmentation and summarization. In *IEEE International Conference on Multimedia and Expo*, volume 3, pages 1329–1332, New York, NY, July 2000.
- [22] A. Dailianas, R.B. Allen, and P. England. Comparison of automatic video segmentation algorithms. In *SPIE Conference on Integration Issues in Large Commercial Media Delivery Systems*, volume 2615, pages 2–16, Philadelphia, PA, October 1995.
- [23] C.H. Demarty and S. Beucher. Morphological tools for indexing video documents. In *IEEE International Conference on Multimedia Computing and Systems*, volume 2, pages 991–992, Florence, Italy, June 1999.
- [24] M.S. Drew, J. Wei, and Z.N. Li. Illumination-invariant image retrieval and video segmentation. *Pattern Recognition*, 32(8):1369–1388, August 1999.
- [25] R. Dugad, K. Ratakonda, and N. Ahuja. Robust video shot change detection. In *IEEE Workshop on Multimedia Signal Processing*, pages 376–381, Redondo Beach, CA, December 1998.
- [26] S. Eickeler and S. Müller. Content-based video indexing of tv broadcast news using hidden markov models. In *International Conference on Acoustics, Speech, and Signal Processing*, pages 2997–3000, Phoenix, AZ, March 1999.
- [27] S. Eickeler and G. Rigoll. A novel error measure for the evaluation of video indexing systems. In *International Conference on Acoustics, Speech, and Signal Processing*, Istanbul, Turkey, June 2000.
- [28] O. Fatemi, S. Zhang, and S. Panchanathan. Optical flow based model for scene cut detection. In *Canadian Conference on Electrical and Computer Engineering*, volume 1, pages 470–473, Calgary, Canada, May 1996.
- [29] A.M. Ferman and A.M. Tekalp. Efficient filtering and clustering for temporal video segmentation and visual summarization. *Journal of Visual Communication and Image Representation*, 9(4):336–351, December 1998.
- [30] W.A.C. Fernando, C.N. Canagarajah, and D.R. Bull. Fade and dissolve detection in uncompressed and compressed video sequences. In *IEEE International Conference on Image Processing*, page 27AP2, Kobe, Japan, October 1999.
- [31] U. Gargi and R. Kasturi. An evaluation of color histogram-based methods in video indexing. In *International Workshop on Image Databases and Multimedia Search*, pages 75–82, Amsterdam, The Netherlands, August 1996.

- [32] U. Gargi, R. Kasturi, and S.H. Strayer. Performance characterization of video-shot-change detection methods. *IEEE Transactions on Circuits and Systems for Video Technology*, 10(1):1–13, February 2000.
- [33] J.M. Gauch, S. Gauch, S. Bouix, and X. Zhu. Real time video scene detection and classification. *Information Processing and Management*, 35(3):401–420, May 1999.
- [34] Y. Gong and X. Liu. Video shot segmentation and classification. In *IAPR International Conference on Pattern Recognition*, volume 1, pages 860–863, Barcelona, Spain, September 2000.
- [35] B. Gunsel, A.M. Ferman, and A.M. Tekalp. Temporal video segmentation using unsupervised clustering and semantic object tracking. *Journal of Electronic Imaging*, 7(3):592–604, July 1998.
- [36] N. Haering, N. da Vitoria Lobo, R. Qian, and I Sezan. A framework for designing event detectors. In *Asian Conference on Computer Vision*, Taipei, Taiwan, January 2000.
- [37] A. Hampapur, R.C. Jain, and T. Weymouth. Production model based digital video segmentation. *International Journal of Multimedia Tools and Applications*, 1(1):9–46, March 1995.
- [38] K.J. Han and A.H. Tewfik. Eigen-image based video segmentation and indexing. In *IEEE International Conference on Image Processing*, pages 538–541, Santa Barbara, CA, October 1997.
- [39] S.H. Han, I.S. Kweon, C.Y. Kim, and Y.S. Seo. Bayesian shot detection using structural weighting. In *IAPR International Workshop on Machine Vision Applications*, pages 95–98, Tokyo, Japan, November 2000.
- [40] A. Hanjalic and H.J. Zhang. Optimal shot boundary detection based on robust statistical models. In *IEEE International Conference on Multimedia Computing and Systems*, volume 2, pages 710–714, Florence, Italy, June 1999.
- [41] W.J. Heng and K.N. Ngan. Integrated shot boundary detection using object-based technique. In *IEEE International Conference on Image Processing*, volume 3, pages 289–293, Kobe, Japan, October 1999.
- [42] F. Idris and S. Panchanathan. Review of image and video indexing techniques. *Journal of Visual Communication and Image Representation*, 8(2):146–166, June 1997.
- [43] O. Javed, S. Khan, Z. Rasheed, and M. Shah. A framework for segmentation of interview videos. In *IASTED International Conference on Internet and Multimedia Systems and Applications*, Las Vegas, CA, November 2000.
- [44] H. Jiang, A.S. Helal, A.K. Elmagarmid, and A. Joshi. Scene change detection techniques for video database systems. *Multimedia Systems*, 6(3):186–195, May 1998.



- [45] R. Kasturi and R.C. Jain. Dynamic vision. In R. Kasturi and R.C. Jain, editors, *Computer Vision: Principles*, pages 469–480. IEEE Computer Society Press, Washington, 1991.
- [46] S. Kim and R.H. Park. A novel approach to scene change detection using a cross entropy. In *IEEE International Conference on Image Processing*, volume 3, pages 937–940, Vancouver, Canada, September 2000.
- [47] W.X. Kong, X.F. Ding, H.Q. Lu, and S.D. Ma. Improvement of shot detection using illumination invariant metric and dynamic threshold selection. In *International Conference on Visual Information Systems*, pages 277–282, Amsterdam, The Netherlands, June 1999.
- [48] I. Koprinska and S. Carrato. Temporal video segmentation: A survey. *Signal Processing: Image Communication*, 16(5):477–500, January 2001.
- [49] S. Lawrence, D. Ziou, M.F. Auclair-Fortier, and S. Wang. Motion insensitive detection of cuts and gradual transitions in digital videos. to appear in *Pattern Recognition Letters*, 2002.
- [50] C.M. Lee and M.C. Ip. A robust approach for camera break detection in color video sequences. In *IAPR International Workshop on Machine Vision Applications*, pages 502–505, Kawasaki, Japan, December 1994.
- [51] H.C. Lee, C.W. Lee, and S.D. Kim. Abrupt shot change detection using an unsupervised clustering of multiple features. In *International Conference on Acoustics, Speech, and Signal Processing*, volume 3, pages 2015–2018, Istanbul, Turkey, June 2000.
- [52] M.S. Lee, Y.M. Yang, and S.W. Lee. Automatic video parsing using shot boundary detection and camera operation analysis. *Pattern Recognition*, 34(3):711–719, March 2001.
- [53] S. Lefèvre, J. Holler, and N. Vincent. Real time temporal segmentation of compressed and uncompressed dynamic colour image sequences. In *International Workshop on Real Time Image Sequence Analysis*, pages 56–62, Oulu, Finland, August 2000.
- [54] D. Li and H. Lu. Efficient scene change detection through stdd (single threshold double direction). In *IASTED International Conference on Modelling and Simulation*, Pittsburgh, PA, May 2000.
- [55] Z.N. Li and J. Wei. Spatio-temporal joint probability images for video segmentation. In *IEEE International Conference on Image Processing*, volume 2, pages 295–298, Vancouver, Canada, September 2000.
- [56] R. Lienhart. Comparison of automatic shot boundary detection algorithms. In *SPIE Conference on Storage and Retrieval for Image and Video Databases*, volume 3656, pages 290–301, San Jose, CA, January 1999.

- [57] R. Lienhart. Reliable transition detection in videos: A survey and practitioner's guide. *International Journal of Image and Graphics*, 1(3):469–486, July 2001.
- [58] T. Liu, X. Zhang, D. Wang, J. Feng, and K.T. Lo. Inertia-based cut detection technique: A step to the integration of video coding and content-based retrieval. In *IEEE International Conference on Signal Processing*, pages 1018–1025, Beijing, China, August 2000.
- [59] M.K. Mandal, F. Idris, and S. Panchanathan. A critical evaluation of image and video indexing techniques in the compressed domain. *Image and Vision Computing*, 17(7):513–529, May 1999.
- [60] S. Mann and R.W. Picard. Video orbits of the projective group: A simple approach to featureless estimation of parameters. *IEEE Transactions on Image Processing*, 6(9):1281–1295, September 1997.
- [61] A. Nagasaka and Y. Tanaka. Automatic video indexing and full-video search for object appearances. In *IFIP Working Conference on Visual Database Systems*, pages 113–127, Budapest, Hungary, October 1991.
- [62] J. Nam and A.H. Tewfik. Combined audio and visual streams analysis for video sequence segmentation. In *International Conference on Acoustics, Speech, and Signal Processing*, volume 4, pages 2665–2668, Munich, Germany, April 1997.
- [63] J. Nam and A.H. Tewfik. Wipe transition detection using polynomial interpolation. In *SPIE Conference on Storage and Retrieval for Media Databases*, volume 4315, pages 231–241, San Jose, CA, January 2001.
- [64] M.R. Naphade, P. Mehrotra, A.M. Ferman, J. Warnick, T.S. Huang, and A.M. Tekalp. A high-performance shot boundary detection algorithm using multiple cues. In *IEEE International Conference on Image Processing*, volume 2, pages 884–887, Chicago, IL, October 1998.
- [65] J.H. Oh, K.A. Hua, and N. Liang. A content-based scene change detection and classification technique using background tracking. In *SPIE Conference on Multimedia Computing and Networking*, volume 3969, pages 254–265, San Jose, CA, January 2000.
- [66] C. O'Toole, A. Smeaton, N. Murphy, and S. Marlow. Evaluation of automatic shot boundary detection on a large video test suite. In *Conference on Challenge of Information Retrieval*, Newcastle, UK, February 1999.
- [67] K. Otsuji and Y. Tonomura. Projection-detecting filter for video cut detection. *Multimedia Systems*, 1(5):205–210, March 1994.
- [68] S.V. Porter, M. Mirmehdi, and B.T. Thomas. Video cut detection using frequency domain correlation. In *IAPR International Conference on Pattern Recognition*, volume 3, pages 413–416, Barcelona, Spain, September 2000.

- [69] D. Pye, N.J. Hollinghurst, T.J. Mills, and K.R. Wood. Audio-visual segmentation for content-based retrieval. In *International Conference on Spoken Language Processing*, Sydney, Australia, December 1998.
- [70] G. Quénot and P. Mulhem. Two systems for temporal video segmentation. In *European Workshop on Content Based Multimedia Indexing*, pages 187–194, Toulouse, France, October 1999.
- [71] W. Ren, M. Sharma, and S. Singh. Automated video segmentation. In *International Conference on Information, Communication, and Signal Processing*, Singapore, October 2001.
- [72] R. Ruiloba, P. Joly, S. Marchand-Maillet, and G. Quénot. Towards a standard protocol for the evaluation of video-to-shots segmentation algorithms. In *European Workshop on Content Based Multimedia Indexing*, pages 41–48, Toulouse, France, October 1999.
- [73] B. Sabata and M. Goldszmidt. Fusion of multiple cues for video segmentation. In *International Conference on Information Fusion*, Sunnyvale, CA, July 1999.
- [74] J.M. Sanchez, X. Binefa, J. Vitria, and P. Radeva. Local color analysis for scene break detection applied to tv commercials recognition. In *International Conference on Visual Information Systems*, pages 237–244, Amsterdam, The Netherlands, June 1999.
- [75] B. Shahraray. Scene change detection and content-based sampling of video sequences. In *SPIE Conference on Digital Video Compression: Algorithms and Technologies*, volume 2419, pages 2–13, San Jose, CA, February 1995.
- [76] R. Silva, J. Gomes, and L. Velho. Segmentation of video sequences using volumetric image processing. In *Eurographics Workshop on Multimedia*, Milan, Italy, September 1999.
- [77] A.F. Smeaton, G. Gormley J. Gilvarry and, B. Tobin, S. Marlow, and N. Murphy. An evaluation of alternative techniques for automatic detection of shot boundaries in digital video. In *Irish Machine Vision and Image Processing Conference*, pages 45–62, Dublin, Ireland, September 1999.
- [78] D. Swanberg, C.F. Shu, and R. Jain. Knowledge guided parsing and retrieval in video databases. In *SPIE Conference on Storage and Retrieval for Image and Video Databases*, volume 1908, pages 13–24, San Jose, CA, February 1993.
- [79] Y. Taniguchi, A. Akutsu, and Y. Tonomura. Panorama excerpts: Extracting and packing panoramas for video browsing. In *ACM International Conference on Multimedia*, pages 427–436, Seattle, WA, November 1997.
- [80] Y. Tonomura and S. Abe. Content oriented visual interface using video icons for visual database systems. *Journal of Visual Languages and Computing*, 1(2):183–198, June 1990.

- [81] B.T Truong, C. Dorai, and S. Venkatesh. New enhancements to cut, fade, and dissolve detection processes in video segmentation. In *ACM International Conference on Multimedia*, pages 219–227, Los Angeles, CA, October 2000.
- [82] H. Ueda, T. Miyatake, and S. Yoshizawa. Impact: An interactive natural-motion-picture dedicated multimedia authoring system. In *ACM Conference on Human Factors in Computing Systems*, pages 343–350, New Orleans, LO, April 1991.
- [83] N. Vasconcelos and A. Lippman. Statistical models of video structure for content analysis and characterization. *IEEE Transactions on Image Processing*, 9(1):3–19, January 2001.
- [84] A. Vellakal and C.C.J. Kuo. Joint spatial-spectral indexing for image retrieval. In *IEEE International Conference on Image Processing*, volume 3, pages 867–870, Lausanne, Switzerland, September 1996.
- [85] W. Xiong and J.C.M. Lee. Efficient scene change detection and camera motion annotation for video classification. *Journal of Computer Vision and Image Understanding*, 71(2):166–181, August 1998.
- [86] A. Yilmaz and M. Shah. Shot detection using principle coordinate system. In *IASTED International Conference on Internet and Multimedia Systems and Applications*, Las Vegas, CA, November 2000.
- [87] H. Yu, G. Bozdagi, and S. Harrington. Feature-based hierarchical video segmentation. In *IEEE International Conference on Image Processing*, volume 2, pages 498–501, Santa Barbara, CA, October 1997.
- [88] Y. Yusoff, W. Christmas, and J. Kittler. Video shot cut detection using adaptative thresholding. In *British Machine Vision Conference*, pages 362–371, Bristol, UK, September 2000.
- [89] Y. Yusoff, J. Kittler, and W. Christmas. Combining multiple experts for classifying shot changes in video sequences. In *IEEE International Conference on Multimedia Computing and Systems*, volume 2, pages 700–704, Florence, Italy, June 1999.
- [90] R. Zabih, J. Miller, and K. Mai. A feature-based algorithm for detecting and classifying production effects. *Multimedia Systems*, 7(2):119–128, March 1999.
- [91] H.J. Zhang, A. Kankanhalli, and S.W. Smoliar. Automatic partitioning of full-motion video. *Multimedia Systems*, 1(1):10–28, July 1993.
- [92] W. Zhao, J. Wang, D. Bhat, K. Sakiewicz, N. Nandhakumar, and W. Chang. Improving color based video shot detection. In *IEEE International Conference on Multimedia Computing and Systems*, volume 2, pages 752–756, Florence, Italy, June 1999.
- [93] D. Zugaj and P. Bouthemy. Wipe detection in the temporal segmentation of video. In *European Workshop on Content Based Multimedia Indexing*, Toulouse, France, October 1999.

Table 1: Complexity and relative computation time of pixel-based methods

Method and reference(s)	Complexity	Time
Simple interframe difference [61]	$O(\mathcal{P})$	1.00
Template matching [61]	$O(3\mathcal{P})$	2.11
Colour template matching [61]	$O(9\mathcal{P})^*$	5.29
Boolean difference [91]	$O(2\mathcal{P})$	1.80
Improved boolean difference [91]	$O(4\mathcal{P})$	2.60
Normalized difference energy [71]	$O(5\mathcal{P})$	3.00
Normalized sum of absolute differences [71]	$O(5\mathcal{P})$	2.58
Pixel labelling [79]	$O(15\mathcal{P})$	
Evolution of temporal derivative of pixel intensities [49]	$O(8\mathcal{P})_+$	

Table 2: Complexity and relative computation time of histogram-based methods

Method and reference(s)	Complexity	Time
Histogram creation	$O(\mathcal{P})$	
Color histogram creation	$O(3\mathcal{P})$	
Histogram difference [80, 61]	$O(3\mathcal{N})$	4.29
Chrominance histogram difference [31]	$O(6\mathcal{N})^*$	55.62
Highest histogram difference [69]	$O(9\mathcal{N})^*$	16.47
Hue histogram difference [5]	$O(4\mathcal{N})$	111.56
Normalized histogram difference [5]	$O(4\mathcal{N})$	12.18
Cosine similarity measure [66, 18]	$O(4\mathcal{N})$	45.84
Genetic algorithm [21] (colour histogram only)	$O(9\mathcal{N})_+$	
Twin comparison [91]	$O(3\mathcal{N})$	
Squared histogram difference [71]	$O(6\mathcal{N})$	11.18
Kolmogorov-Smirnov statistic [71]	$O(3\mathcal{N})$	4.29
Ponderation relative to colour importance [22]	$O(3\mathcal{P} + 9\mathcal{N})^*$	19.58
Ponderation after learning step [92]	$O(15\mathcal{N})_+$	
Average colour [31]	$O(6\mathcal{N})^*$	21.56
Reference colour table [31]	$O(8\mathcal{N})$	12.29
Histogram intersection [31, 36, 43]	$O(2\mathcal{N})$	4.29
Normalized histogram intersection [31]	$O(4\mathcal{N})$	12.18
Original $\chi^2$ test [61]	$O(5\mathcal{N})$	12.29
$\chi^2$ test normalized by max [22]	$O(5\mathcal{N})$	12.27
$\chi^2$ test on colour histograms [35]	$O(15\mathcal{N})^*$	
Cross entropy [71, 46]	$O(4\mathcal{N})$	12.29
Divergence [71, 46]	$O(8\mathcal{N})$	12.27
Kullback Liebler distance [71, 46]	$O(11\mathcal{N})$	12.29
Bhattacharya distance [71]	$O(3\mathcal{N})$	12.27

Table 3: Complexity and relative computation time of block-based methods

Method and reference(s)	Complexity	Time
Similarity ratio [45]	$O(3\mathcal{P} + 15\mathcal{B})$	85.29
Yakimovsky likelihood ratio [71]	$O(3\mathcal{P} + 13\mathcal{B})$	88.67
Freund statistic	$O(3\mathcal{P} + 11\mathcal{B})$	79.20
Block differences in the HSV colour space [52]	$O(2\mathcal{P} + 10\mathcal{B})^*$	179.05
Colour histogram comparison [78]	$O(15\mathcal{N}\mathcal{B}) + 3\mathcal{B})^*$	206.42
Selective RGB histogram comparison [61]	$O(4\mathcal{N}\mathcal{B} + 0, 5\mathcal{B})_+$	208.42
Rate of correlation change [82]	$O(4\mathcal{N}\mathcal{B} + 2\mathcal{B})$	
Block-based normalized histogram difference [5]	$O(4\mathcal{N}\mathcal{B} + 2\mathcal{B})$	233.38
Selective HSV histogram comparison [50]	$O(6\mathcal{N}\mathcal{B} + 2\mathcal{B} + 4\mathcal{P})^*$	300.00
Block-based HSI histogram comparison [13]	$O(4\mathcal{N}\mathcal{B} + \mathcal{B})^*$	
L*u*v* histogram intersection [19]	$O(12\mathcal{N}\mathcal{B} + 2\mathcal{B})^*$	
Histogram difference and likelihood ratio [25]	$O(3\mathcal{P} + 3\mathcal{N} + 21\mathcal{B})$	
Neighbourhood colour ratio [2]	$O(5\mathcal{P} + 8\mathcal{B})$	45.93
RGB block dissimilarity [23]	$O(9\mathcal{P} + 3\mathcal{B})^*$	
HSV block dissimilarity [53]	$O(2\mathcal{P} + 8\mathcal{B})^*$	114.25
Temporal and spatial subsampling [85]	$O(\mathcal{P} + 5\mathcal{B})$	36.45

Table 4: Complexity and relative computation time of feature-based methods

Method and reference(s)	Complexity	Time
Moment invariants [12]	$O(23\mathcal{P})$	
Edge change ratio [90, 77]	$O(26\mathcal{P})_+$	240.67
Edge-based contrast [56]	$O(25\mathcal{P})$	
Histogram difference, edge points count, double chromatic difference [87]	$O(27\mathcal{P} + 3\mathcal{N})$	
Planar points [76]	$O(51\mathcal{P})$	
Transition modelling using centroids [37]	$O(10\mathcal{P})$	11.89
Transition modelling using histograms [1]	$O(2\mathcal{N}^2 + 10\mathcal{N})$	
Histogram of pixel-pair differences [6]	$O(8\mathcal{P}_+)$	
Transition modelling using mean and variance [56, 10, 9, 81, 30]	$O(3\mathcal{P})$	
Principal coordinate system [86]	$O(18\mathcal{P})^*$	35.95
Joint probability images [55]	$O(2\mathcal{P} + \mathcal{N}^2 + \mathcal{N})$	
Hidden Markov Model with 3 features [26]	$O(13\mathcal{P} + 11\mathcal{N})$	



Table 5: Complexity and relative computation time of motion-based methods

Method and reference(s)	Complexity	Time
Global motion [20, 16, 93, 60]	$O(20\mathcal{P})_+$	80.80
Motion vectors [8, 75, 58]	$O(75\mathcal{P})_+$	2312.73
Optical flow [28]	$O(105\mathcal{P})_+$	
Frequency domain correlation [68]	$O(75\mathcal{P})_+$	105.65

Table 6: Complexity and relative computation time of other methods

Method and reference(s)	Complexity	Time
$\chi^2$ test and edge analysis [70]	$O(26\mathcal{P} + 5\mathcal{N})$	243.07
Average brightness, pixel value change and colour distribution change [33]	$O(4\mathcal{P} + 3\mathcal{N})$	16.47
Combination of 5 algorithms [89]	$O(85\mathcal{P} + 9\mathcal{N} + 17\mathcal{B})_+$	
K-means using colour histogram and pixel intensity differences [64, 29]	$O(4\mathcal{P} + 9\mathcal{N})^*$	
K-means using gray-level histogram and pixel intensity differences [51]	$O(3\mathcal{P} + 3\mathcal{N})$	

## List of Tables

1	Complexity and relative computation time of pixel-based methods . .	44
2	Complexity and relative computation time of histogram-based methods	45
3	Complexity and relative computation time of block-based methods . .	46
4	Complexity and relative computation time of feature-based methods .	47
5	Complexity and relative computation time of motion-based methods .	48
6	Complexity and relative computation time of other methods . . . . .	49