



**[biblio.ugent.be](http://biblio.ugent.be)**

The UGent Institutional Repository is the electronic archiving and dissemination platform for all UGent research publications. Ghent University has implemented a mandate stipulating that all academic publications of UGent researchers should be deposited and archived in this repository. Except for items where current copyright restrictions apply, these papers are available in Open Access.

This item is the archived peer-reviewed author-version of:

Compressed-Domain Shot Boundary Detection for H.264/AVC Using Intra Partitioning Maps

Sarah De Bruyne, Jan De Cock, Chris Poppe, Charles-Frederik Hollemeersch, Peter Lambert and Rik Van de Walle

In: Lecture Notes in Computer Science, 6523(2011), 29-39, 2011.

Optional: <http://www.springerlink.com/content/v0h85302q6761g07/>

**To refer to or to cite this work, please use the citation to the published version:**

**Sarah De Bruyne, Jan De Cock, Chris Poppe, Charles-Frederik Hollemeersch, Peter Lambert and Rik Van de Walle (2011). Compressed-Domain Shot Boundary Detection for H.264/AVC Using Intra Partitioning Maps. *Lecture Notes in Computer Science* 6523(2011) 29-39. DOI: 10.1007/978-3-642-17832-0\_4**

# Compressed-domain shot boundary detection for H.264/AVC using intra partitioning maps

Sarah De Bruyne, Jan De Cock, Chris Poppe, Charles-Frederik Hollemeersch,  
Peter Lambert, and Rik Van de Walle

Ghent University - IBBT,  
Department of Electronics and Information Systems - Multimedia Lab  
Gaston Crommenlaan 8 bus 201, B-9050 Ledeborg-Ghent, Belgium  
[sarah.debruyne@ugent.be](mailto:sarah.debruyne@ugent.be)  
<http://multimedialab.elis.ugent.be>

**Abstract.** In this paper, a novel technique for shot boundary detection operating on H.264/AVC-compressed sequences is presented. Due to new and improved coding tools in H.264/AVC, the characteristics of the obtained sequences differ from former video coding standards. Although several algorithms working on this new standard are already proposed, the presence of IDR frames can still lead to a low accuracy for abrupt transitions. To solve this issue, we present the motion-compensated intra partitioning map which relies on the intra partitioning modes and the motion vectors present in the compressed video stream. Experimental results show that this motion-compensated map achieves a high accuracy and exceeds related work.

**Keywords:** Shot boundary detection, video analysis, compressed domain, H.264/AVC

## 1 Introduction

During the last decades, a significant increase in the use and availability of digital multimedia content can be witnessed. Unfortunately, these video collections often lack information related to the structure and the actual content of the video. When accessing these video streams in case no metadata is available, time-consuming, sequential scanning is the only option. As a consequence, to facilitate multimedia consumption, intensive research has been done in the domain of indexing, retrieval, browsing, and summarization. Since the identification of the temporal structure of video is an essential task for many video indexing and retrieval applications, the first step commonly taken for video analysis is shot boundary detection as shots are the basic units for a large majority of video content analysis algorithms [5]. According to whether the transition between consecutive shots is abrupt or not, boundaries are classified as cuts or gradual transitions.

In order to preserve storage space and to reduce bandwidth constraints, most video data is available in compressed form. By relying on compressed-domain features which can be extracted directly from the compressed bitstream,

time-consuming decompression can be avoided and coarse but potentially useful information present in the bitstream can efficiently be reused. Consequently, compressed-domain algorithms for shot boundary detection are gaining importance. In the past, many compressed-domain algorithms were proposed which rely on the MPEG-1 Video and MPEG-2 Video standards. However, as the H.264/AVC video coding standard [12] performs significantly better than any prior standard in terms of coding efficiency, more video content will be coded in this video format in the future. Its superior compression performance can mainly be attributed to the new or improved coding tools. However, these coding tools influence the compressed domain features to a great extent and render prior algorithms working on MPEG-1 Video and MPEG-2 Video obsolete. As a consequence, recently, efforts have been undertaken to design new shot boundary detection algorithms working on H.264/AVC.

The outline of this paper is as follows. Section 2 addresses related work and remaining issues in the area of shot boundary detection algorithms operating on H.264/AVC-compressed video streams. Section 3 introduces our novel algorithm to detect shot boundaries, whereas results are provided in Section 4. Conclusions are drawn in Section 5.

## 2 Related work

### 2.1 General techniques

In literature, most of the algorithms work on MPEG-1 Video and MPEG-2 Video. On the one hand, DCT coefficients are exploited. Arman *et al.* [1] compare a subset of DCT coefficients of two successive I frames to calculate the frame differences as the coefficients in the frequency domain are mathematically related to the spatial domain. The temporal resolution of this type of techniques is low, resulting in an increased amount of false alarms when camera motion is present. Furthermore, Yeo and Liu defined the concept of DC images [13], which are spatially reduced versions of the original image and which are generated by only taking into account the first DCT coefficient in each block, i.e., the DC coefficient, and motion vectors. Based on these DC images, similarity metrics defined for color features in the pixel domain can be modified to operate in the compressed domain.

On the other hand, the distribution of the different macroblock types and motion information [4, 11] can also be used as features to detect shot boundaries. When an abrupt transition occurs between two successive P pictures, it is expected that a significant amount of macroblocks in the second frame is intra coded since these macroblocks cannot be predicted well from prior reference frames. The prediction directions of motion vectors in intermediate B frames can then be utilized to detect the exact location of the transition.

Although most algorithms mainly focus on the detection of abrupt transitions, the aforementioned features can also be used to detect gradual transitions [2]. Due to the large variety in terms of effects and duration, the accuracy for gradual transitions is typically inferior to the abrupt transitions.

## 2.2 Algorithms for H.264/AVC

H.264/AVC contains a number of new or improved coding tools which have a major impact on the aforementioned shot boundary detection algorithms. Firstly, intra prediction in H.264/AVC is conducted in the spatial domain by relying on neighboring samples of previously-decoded blocks in order to reduce the spatial redundancy in images. As such, DC coefficients in intra-coded pictures no longer represent average energy, but only represent an energy difference. Consequently, shot boundary detection algorithms working on DC images can no longer be applied to H.264/AVC bitstreams. The second feature, multiple reference picture motion compensation, allows an encoder to not only use the previous and following reference frame during encoding, but makes it possible to also use additional priorly-decoded frames as reference. As this results in vagueness about random access in the bitstream, the concept of Instantaneous Decoding Refresh (IDR) was introduced [12]. This special I frame indicates that no subsequent pictures in the bitstream will require references to pictures prior to the IDR picture in decoding order. The prediction chain is broken; hence it is insufficient to only rely on reference directions to detect shot boundaries.

Up to now, a few algorithms working on H.264/AVC-compressed bitstreams have been published. In [8], Kim *et al.* define a dissimilarity metric for I frames by relying on macroblock partitions (i.e., Intra\_4×4 and Intra\_16×16 prediction). A more complex discontinuity metric based on solely I frames is proposed by Kuo and Lo where intra mode histogram distances are calculated based on the different intra prediction mode directions [9]. However, the exact location of a shot boundary cannot be determined by these two algorithms as information from P and B frames is not taken into consideration. Extensions on this second algorithm are proposed in [10] and [14] to deal with all type of frames. Firstly, these algorithms exploit the intra prediction histograms to locate potential groups of pictures (GOPs) where the probability of a shot boundary is high. Secondly, the different inter prediction modes and motion vector directions of the intermediate P and B frames are used to locate the exact location of the transition by making use of thresholds or Hidden Markov Models. However, due to the presence of IDR frames, shot boundaries located just before the IDR frames are falsely detected in case the dissimilarity between the two I frames is relatively large. Furthermore, when the dissimilarity metric is low but both I frames actually belong to different shots, the temporal prediction chain is not considered, leading to missed detections.

In our previous work [3], these two problems were tackled. To locate the abrupt transitions, the temporal dependencies between successive frames are examined first. Secondly, when encountering an IDR frame breaking this prediction chain, spatial dissimilarities are considered. In contrast to the aforementioned related algorithms, the changing characteristics of the content are taken into account by constructing a “static intra partitioning map”. In particular, the static intra partitioning map is updated when new intra-coded macroblocks residing at intermediate, inter-coded frames are encountered. Gradual changes are detected by relying on the percentage of intra-coded macroblocks. Since fast global and

local motion can result in similar patterns as gradual transitions, motion-activity intensity is considered as well to identify the exact origin of the content change.

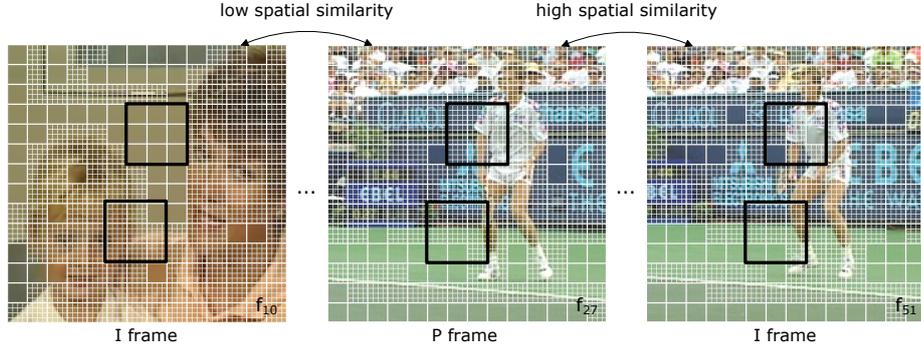
Although the static intra partitioning map already solves several issues, its accuracy is still inferior to video streams compressed without IDR frames. When examining the origin of the false alarms, it can be seen that they mostly occur when the content slightly changes as a result of motion, which is generally compensated by using inter prediction. In the next section, we will extend the algorithm proposed in [3] by introducing the motion-compensated intra partition map which overcomes this problem by relying on motion compensation techniques.

### 3 Intra partitioning maps

Whereas the abrupt transitions located in between two inter-coded frames can be detected by analyzing the main reference directions of the frames, another technique is required to determine whether a shot boundary is present in between an IDR frame and its preceding frames. To solve the issue of the broken temporal prediction chain, algorithms designed for prior video coding standards would typically make use of DC images to compare the spatial characteristics of successive frames. However, due to the introduction of spatial intra prediction in H.264/AVC, these iconic versions of the content can no longer be generated without further decoding the bitstream. Therefore, we employ intra partitioning maps which reflects the spatial dissimilarities between frames based on the selected intra partitioning modes.

H.264/AVC supports multiple intra macroblock partitions, i.e., Intra\_4×4, Intra\_8×8 and Intra\_16×16. The first mode is generally selected by an encoder in case of significant detail, whereas the Intra\_16×16 mode is preferred for smooth areas. When coding high-resolution video sequences using the High profile of H.264/AVC, Intra\_8×8 will also be selected and mainly replace Intra\_4×4 modes to code the high-textured areas. As the subdivision in different macroblock partitions roughly reflects the detail of the content, comparing the distribution of two frames can be used to estimate the spatial dissimilarities, as illustrated in Fig. 1.

Comparing the current I frame with the previous I frame is not recommended as the content can change significantly between these two intra-coded frames. For example, a shot boundary can be located at an intermediate P or B frame, new objects can appear, or camera motion can occur (Fig. 1). Therefore, we introduce an *intra partitioning map*  $M_i$  indicating which intra partitioning modes most likely correspond to the content of the intermediate, inter-predicted frames. This partitioning map can be constructed in different ways. In [3], the static intra partitioning map was introduced, which considers all intra-coded macroblocks in intra- as well as inter-coded frames. In this paper, the motion-compensated intra partitioning map is introduced, which extends the static map by including motion information.



**Fig. 1.** Distribution of Intra\_4×4 and Intra\_16×16 macroblocks. Although the second frame is a P frame, it is mainly intra coded as it is the first frame of a new shot. As such, it is important to update the intra partitioning maps with information from inter-coded frames.

*Static intra partitioning map* A first approach to construct an intra partitioning map for the inter-coded frame  $f_i$  is by remembering for each macroblock position the intra partitioning mode which was last encountered. To put it differently, the intra partitioning map  $M_i$  of frame  $f_i$  is constructed by updating the previous map  $M_{i-1}$  with the partitioning modes of the intra-coded macroblocks in the current frame. As such, this map can be used to represent the spatial distribution of the content of the current frame, in spite of the fact that this frame mainly contains inter-coded macroblocks. By comparing the current I frame  $f_i$  with the static map  $M_{i-1}$ , the spatial dissimilarity between the current and previous frame can be calculated. However, instead of comparing partitioning modes at corresponding positions, a window of several macroblocks is selected for each macroblock. This way, small movement of objects or the camera will lead to fewer false alarms.

*Motion-compensated intra partitioning map* The downside of the static intra partitioning map is its limited support for motion. In particular, when the difference between two I frames belonging to the same shot is large, resulting from the large distance between these two frames or from relatively fast moving objects or the camera, this window will not be able to cover the displacement of the scene. As a result, the amount of falsely detected shot boundaries will increase. To overcome this problem, motion information can be considered in order to construct a motion-compensated intra partitioning map.

Instead of copying the partitioning modes of the previous intra partitioning map, we propose to use the partitioning modes of the reference blocks to which the motion vectors (MVs) of the current frame point to. The favorable aspect of this motion compensation step is shown in Fig. 2(a) and 2(c) and further explained below. The macroblock in  $f_{186}$  which is marked in red is part of the

low-textured sky (Fig. 2(c)) and would typically result in an `Intra_16×16` partitioning mode when intra prediction would be applied, as verified by coding the sequence with I frames only. When copying the partitioning mode of the macroblock located at the same location in the previous partitioning map, marked by the semi-transparent red rectangle in Fig. 2(a), incorrect spatial information belonging to the high-textured crane would be passed on. However, thanks to the motion compensation step, correct spatial information can now be stored in the partitioning map, which corresponds to the block indicated by the full red rectangle in Fig. 2(a). Obviously, for intra-coded macroblocks, the new partitioning modes are used to update the partitioning map (Fig. 2(d)).

The reference block to which the MV of the current frame points to does not necessarily coincide with macroblock or sub-macroblock boundaries. As such, this reference block can overlap with different partitioning modes, as illustrated by the green block in Fig. 2(a). Therefore, it is undesired to store only one partitioning mode for each block. Instead, we propose to store the likelihood of each possible partitioning mode. As multiple partitioning modes for inter-coded macroblocks exist and each partition contains its own MVs, it is desired to divide the intra partitioning map into a uniform field of basic units of  $4 \times 4$  pixels, which corresponds to the smallest partition for which MVs can change. This subdivision also results in a finer granularity, improving the accuracy of the likelihoods of each partitioning mode.

In order to calculate the likelihoods of the different partitioning modes for an inter-coded block  $b$  in the motion-compensated intra partitioning map  $M_i$ , i.e.,  $M_i[b, mode]$ , the likelihoods of the blocks in the reference map  $M_{ref}$  which overlap with the reference block are considered. To incorporate the percentage of overlap between the motion-compensated reference block and the overlapping block  $r$ , a new variable  $OverlappingSize_r$  is introduced. This leads to the following formula where  $OverlappingBlocks_b$  is defined as the set of overlapping blocks connected to block  $b$ :

$$M_i[b, mode] = \sum_{r \in OverlappingBlocks_b} OverlappingSize_r \cdot M_{ref}[r, mode]. \quad (1)$$

When an intra-coded macroblock is encountered, the sixteen corresponding blocks in the intra partitioning map are set to one for the encountered partitioning mode, whereas the other modes are set to zero.

Let  $f_i$  denote the current I frame and  $M_i$  the intra partitioning map containing the new intra partitioning modes,  $M_{i-1}$  the map of the previous frame, and  $B_4$  the amount of  $4 \times 4$  blocks in a frame. Define  $Modes$  as the set of possible macroblock partitions ( $Modes = \{Intra\_4 \times 4, Intra\_8 \times 8, Intra\_16 \times 16\}$ ). Furthermore, let  $M$  be a placeholder for  $M_i$  and  $M_{i-1}$ , and  $n$  and  $m$  macroblocks. The dissimilarity metric  $\Omega$  used to compare the current I frame  $f_i$  and the preceding motion-compensated intra partitioning map can then be defined by the average of the dissimilarity values of all blocks  $b$  in  $f_i$ .

$$\Omega(f_i) = \frac{1}{B_4} \sum_{b \in f_i} \omega(f_i, b). \quad (2)$$









(a) I frame ( $f_{184}$ )                      (b) next inter-coded reference frame ( $f_{186}$ )

(c) inter-coded macroblocks and MVs in  $f_{186}$                       (d) intra-coded macroblocks in  $f_{186}$

**Fig. 2.** Between the two reference frames  $f_{184}$  and  $f_{186}$  of the Foreman sequence, the camera is panning to the right side. To update the motion-compensated intra partitioning map  $M_{186}$ , the MVs of inter-coded macroblocks in  $f_{186}$  are used to locate the corresponding intra partitioning modes in  $f_{184}$  (a and c). The intra-coded macroblocks in  $f_{186}$  are used to update the intra partitioning modes (d).