

SEGMENTATION AND CLASSIFICATION OF MOVING OBJECTS FOR INTELLIGENT VIDEO SURVEILLANCE

B.SHRAVAN KUMAR,

Research Scholar, JNTUH,India,

Dr.V.USHA SHREE

Professor ECE Dept.

Abstract: Moving object segmentation and classification from compressed video plays an important role for intelligent video surveillance. Compared with H.264/AVC, HEVC introduces a host of new coding features which can be further exploited for moving object segmentation and classification. In this paper, we present a real-time approach to segment and classify moving object using unique features directly extracted from the HEVC compressed domain for video surveillance. In the proposed method, firstly, motion vector interpolation for intra-coded prediction unit and MV outlier removal are employed for preprocessing. Secondly, blocks with non-zero motion vectors are clustered into the connected foreground regions using the four connectivity component labeling algorithm. Thirdly, object region tracking based on temporal consistency is applied to the connected foreground regions to remove the noise regions. The boundary of moving object region is further refined by the coding unit size and prediction unit size. Finally, a person-vehicle classification model using bag of spatial-temporal HEVC syntax words is trained to classify the moving objects, either persons or vehicles. The experimental results demonstrate that the proposed method provides the solid performance and can classify moving persons and vehicles accurately.

Keywords: Compression Domain, Object Segmentation, Object Classification, HEVC, Video Surveillance.

I. INTRODUCTION

Moving object segmentation and classification from video data is one of the most important tasks for intelligent video surveillance. Most computer vision methods for moving object detection and classification assume that the original video frames are available and extract descriptions or features from pixel domain [1-3]. Note that most video content are received or stored in compressed formats encoded with international video coding standards, such as MPEG-2 [4], H.264/AVC [5] and HEVC [6]. To obtain the original video frame, we have to perform video decoding. In video analysis at large scales, such as content analysis and search for a large surveillance network, the complexity of video decoding becomes a major bottleneck of the real-time system. To address this issue, compression-domain approaches have been explored for video content analysis which extracts features directly from the bit stream syntax, such as motion vectors and block coding modes [7]. The major advantage of compression-domain approaches is their low computational complexity since the full-scale decoding and reconstruction of pixels are avoided. Therefore, compressed domain methods are desired for real-time video analysis applications. In this paper, we focus on moving object detection and classification from HEVC compressed surveillance videos. Specifically, by extracting features from

HEVC compressed surveillance video bitstream, the moving objects are located and classified, such as persons or vehicles.

A. Related Work

Recently, a number of moving object segmentation and classification algorithms using motion vector (MV) information in H.264/AVC compression domain have been reported. R. V. Babuet al. [8] introduce a method to accumulate MV information over time for moving object segmentation. Temporally accumulated MVs are further interpolated spatially to obtain a dense field, and the expectation maximization procedure is then applied on the dense motion field for final segmentation. S. D. Bruyneet al. [9] develop a method to analyze the reliability of MVs in H.264/AVC domain. This reliability information along with the MV magnitude is used to segment foreground objects from the background. In [10] and [11], MVs are classified into multiple types, such as background, edge, foreground, and noise. Then, the MVs and their associated class information are used to segment each block. In [12] and [13], global motion is first removed from the motion vector field, and moving object segmentation process is performed on the compensated motion vector field. Y. Chen et al. [14] develop a method to extract moving object regions from compressed domain by using global motion estimation and Markov random field classification. It should be noted that MVs extracted from the compressed bitstream are determined in terms of rate-distortion and they may not represent the true object motion. Therefore, it is difficult to find real moving objects solely based on the MV fields.

To address this issue, other information, such as DCT coefficients and macro-block partitions, are used to detect and track moving objects. C. Poppe et al. [15] propose to use the number of bits consumed by each 4x4 block to detect the moving objects from H.264/AVC videos. F. Porikiet al. [16] present a segmentation method that takes advantage of the inter-frame motion and intra-frame DCT coefficients embedded in MPEG videos. M. Laumer et al. [17] propose to use (sub-)block types to refine the level of motion vector for each block. P. Dong et al. [18] develop a moving object segmentation and tracking method by adaptively using the information from motion vectors, DCT coefficients and prediction modes. H. Sabirin et al. [19] propose a spatiotemporal graph-based method for detecting and tracking moving objects by treating the encoded blocks with non-zero motion vectors and/or non-zero residues in H.264/AVC bit streams as potential parts of moving objects. S. H. Khatoonabadi et al. [20] present a method to track moving objects in H.264/AVC compressed video using a spatiotemporal Markov random field (ST-MRF) model, which naturally integrates spatial and temporal aspects of the object motion using motion vectors and block

coding modes. As the newest international standard for video coding, HEVC provides equivalent subjective quality with about 50% bitrate reduction compared to H.264/AVC high profile [21]. HEVC has adopted a host of new encoding features and tools, such as coding unit and prediction unit, which can be exploited for moving object segmentation and classification in compressed domain. However, little work has been done on moving object analysis directly from HEVC compressed videos.

H. Li et al. [22] present a rapid abnormal event detection method for video surveillance by using the syntax features extracted from HEVC compressed domain. Both the syntax information extracted from the HEVC compressed domain and the color information in pixel domain are used to segment the foreground and background region [23, 31, 32]. D. Park et al. [24] propose to extend the ST-MRF model from H.264/AVC compressed domain to HEVC compressed domain for moving object tracking. M. Moriyama et al. [30] propose to conduct this segmentation procedure after the temporal sub-sampling of the video sequence. However, the unique features, such as coding unit and prediction unit of HEVC, are not explored for moving object segmentation and classification. In this paper, we develop a framework for moving object segmentation and classification by using the motion vectors and associated modes directly extracted from HEVC compressed video. Specifically, we focus on surveillance videos whose cameras are stationary. Compared to existing methods in the literature, our work has the following unique aspects and innovations: (1) the unique features in the HEVC compressed domain, such as coding unit and prediction unit, are employed to refine the moving object boundary; (2) the bag of words representation in the HEVC compressed domain is applied to classify the moving persons and vehicles.

B. Overview of the Proposed Method

The overall framework of our system is illustrated in Fig. 1. It consists of two stages: moving object segmentation and person-vehicle classification. For moving object segmentation, firstly, MV interpolation for intra-coded prediction unit (PU) and MV outlier removal are employed for preprocessing. Then, blocks with non-zero motion vectors are clustered into the connected foreground regions by using the four-connectivity component labeling algorithm [25]. Finally, object region tracking with temporal consistency is applied to the connected foreground regions to remove the noise regions. The boundary of moving object region is further refined by using the coding unit (CU) and PU sizes of the blocks.

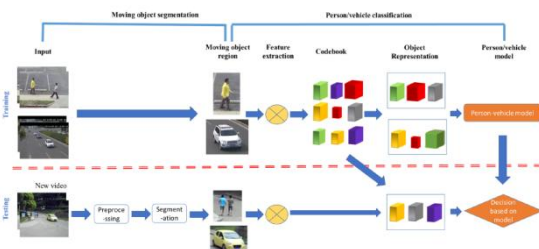


Fig. 1 The proposed framework for HEVC compressed domain moving object segmentation and classification.

For person-vehicle classification, it involves a training phase to learn the person-vehicle model using “bag of spatiotemporal HEVC syntax words” and a testing phase to

apply the learned model to test videos. For the testing phase, we first extract the spatial and temporal information of each 4×4 block to obtain the feature descriptor. Then, the descriptors of all blocks are clustered into a number of codeword. The foreground object is represented by a histogram of the codeword. Finally, for the segmented moving object, we apply the learned person-vehicle model to determine which category to assign. The rest of paper is organized as follows. Section II presents the preprocessing steps. Section III presents the HEVC compressed domain moving object segmentation, tracking and refinement. The compressed domain object classification using bag of HEVC syntax words is described in Section IV. Section V presents the experimental results. Section VI concludes the paper.

II. PREPROCESSING

In HEVC compressed video, one MV is associated with an inter-coded prediction unit (PU). The motion vectors are scaled appropriately to make them independent of the frame type. This is accomplished by dividing the MVs according to the difference between the corresponding frame number and reference frame number (in display order). For example, one MV has values (4,4) for reference frame -1 while another MV in a nearby block has values (8,8) for reference frame -2, these two MV values will be corrected to both be (4,4) after the scaling process. For the PU with two motion vectors, the motion vector with larger length will be selected as the representative motion vector of the PU. In the preprocessing process, the MV interpolation for intra-coded blocks and MV outlier removal are employed before the moving object segmentation and classification.

A. Motion Vector Interpolation for Intra-coded PUs

In order to segment the foreground and background region, it is useful to assign a MV to an intra-coded PU. We propose to select the representative MV from the MVs of its neighboring PUs as its MV. To be specific, the MVs of first-order neighboring PUs (top-left, top, top-right, left, right, bottom-left, bottom, bottom-right) are employed. Fig. 2 shows an example of an intra-coded PU together with its first-order neighboring PUs. In Fig. 2, MVs of all neighboring PUs are collected and stored in MVList. Since one of the neighboring PUs is intracoded, MVList contains seven vectors, which is shown as follows: (1) After MVList is constructed, the next step is to assign a representative MV for the intra-coded PU from MVList. Intracoded PUs usually occur when there is a large motion in the scene. Therefore, we propose to choose the maximum MV from MVList as the MV of the current intra-coded PU. Specifically, when all the first-order neighboring PUs are encoded with intramode, in order to obtain the non-zero MV within the nearby

$$MVList = (MV_1, MV_2, MV_3, MV_4, MV_5, MV_6, MV_7) \quad (1)$$

region, we extend the range of neighborhood to 16×16 (blocks), which is set empirically as being optimal.

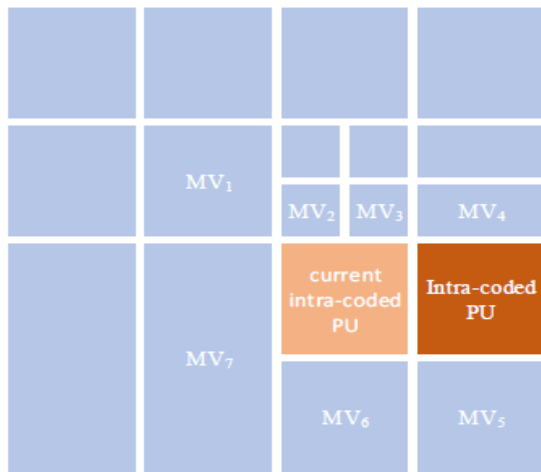


Fig. 2. MV assignment for an intra-coded PU. One of the first-order neighboring PUs is also intra-coded, and the remaining variably sized neighboring PUs have MVs.

B. MV Outlier Removal

Because the MVs from the compressed bitstream are determined in terms of rate-distortion, they may not represent the true object motion. In other words, MV fields may be noisy. In this section, we propose to reduce the motion noise by referring to motion continuity over time and motion coherence within the spatial neighborhood. Three steps are included in the MV outlier removal, which are MV filtering, MV refining, and isolated and small MV removal.

MV Filtering: The original MVs are filtered along the temporal direction. To be specific, the original MVs at the co-located position in the m previous frames and m following frames are employed to filter the original MVs at current frame t . Since the CU and PU sizes at the same position may be different among different frames, MV filtering is operated at 4×4 block level, which is

the minimum size of PU. Let $MV_t^x(k, l)$ and $MV_t^y(k, l)$ represent the original MVs along the horizontal direction and vertical direction for a 4×4 block at position (k, l) at frame t . Then, the filtered MVs: $MV_t^{x'}(k, l)$ and $MV_t^{y'}(k, l)$ can be estimated by

$$MV_t^{x'}(k, l) = \text{floor}\left(\frac{\sum_{i=t-m}^{t+m} MV_i^x(k, l)}{2m+1}\right) \quad (2)$$

$$MV_t^{y'}(k, l) = \text{floor}\left(\frac{\sum_{i=t-m}^{t+m} MV_i^y(k, l)}{2m+1}\right) \quad (3)$$

where $\text{floor}(x)$ represents the operation to get the largest integer less than or equal to x . In the experiment, m is empirically set to 4.

MV Refining: Since the moving objects in the previous frames and following frames have a displacement relative to the moving object in current frame, the filtering process using co-located blocks in the neighboring frames may cause a few non-zero MV noise adjacent to the moving object boundary in current frame. Although the original MVs may be noisy, but if the original MV of current block and most of

its neighboring blocks are both zero, current block has a high probability to belong to background. Figs. 3 shows an example of the original MVs and their associated filtered MVs on frame #53 of the Hall Monitor sequence. In Figs. 3 (a) and (b), the blocks with non-zero MVs and intra mode are marked with green border and red border respectively. As it is shown, the blocks marked with orange circles belong to background. The original MVs of these blocks and most of their neighboring blocks are zero MV. But their associated filtered MVs become non-zero MV.

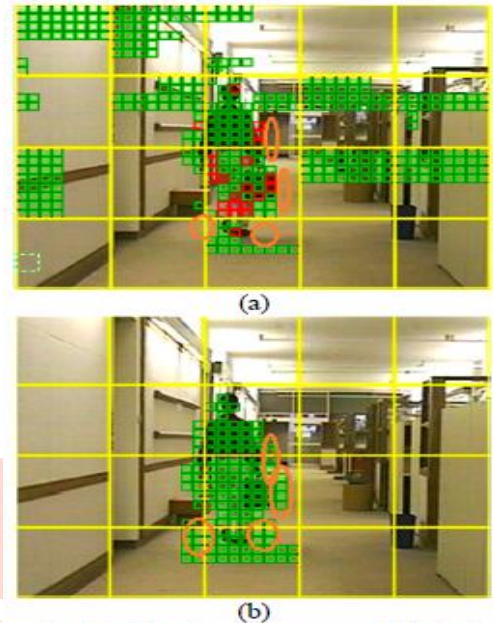


Fig. 3. (a) an example of original motion vectors; (b) their associated filtered motion vectors.

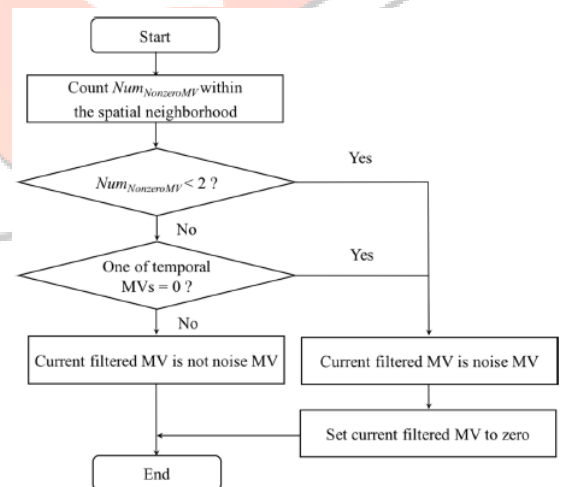


Fig. 4. MV refining for the filtered non-zero MVs.

Therefore, we can use the original MVs to reduce these non-zero MV noise for the filtered MVs. According to the spatial compactness and temporal continuity of the MVs, the original MVs of spatial neighboring PUs and temporal co-located PUs are both involved to make a joint judgment. When the spatial neighboring PU of current PU is out of the frame, the original MV of that PU will be set to zero. The flowchart of the MV refining for filtered non-zero MVs is illustrated in Fig. 4. In Fig. 4, $NumNonzeroMV$ denotes the number of non-zero MVs within the spatial neighborhood, where top PU, bottom PU, left PU and right PU are considered. First, if $NumNonzeroMV$ is less than 2, we assume the condition of spatial compactness is not

satisfied. Second, we check whether the MVs in the co-located temporal PUs from the previous frame and following frame are non-zero or not. If one of the MVs is zero, we assume the condition of temporal continuity is not satisfied. If the condition of the spatial compactness or the condition of temporal continuity is not satisfied, the non-zero filtered MV will be marked as noise motion and set to zero.

Isolated and Small MV Removal: For a foreground moving object, it usually has a connected non-zero MV region and a relatively larger filtered MV, so the PUs with isolated non-zero MVs or smaller MVs have a high probability to be the background PUs. Therefore, we propose to label the PUs with isolated non-zero MVs or small MVs as the background PUs. To be specific, we define one MV as an isolated MV when all the MVs of its spatial neighborhood are zero MV. In addition, we define one MV as a small MV when the MVs of current PU and more than half of its spatial neighboring PU are less than or equal to 1. If one PU is identified as the PU with isolated or small MV, its associated MV will be modified to zero.

III. MOVING OBJECT SEGMENTATION IN HEVC COMPRESSED DOMAIN

After the preprocessing of the MVs, as described in Section II, blocks with non-zero MVs are marked as foreground blocks. These foreground blocks are clustered to the connected foreground regions using the four-connectivity component labeling algorithm [25]. For each foreground region, firstly, we examine its temporal consistency by using object region tracking. Secondly, we refine the boundary of moving object region by using CU and PU sizes of the blocks.

A. Object Region Tracking

In order to examine the temporal consistency of foreground regions, these foreground regions are temporally tracked by using the MVs extracted from HEVC compressed domain. For the i th foreground region O_i^t at frame t , if we find that its corresponding object region continuously describes the same object in backward direction (from t to $t-4$) and in forward direction (from t to $t+4$), then we assume that current foreground region satisfies the condition of temporal consistency. Otherwise, this foreground region will be labeled as the noise region and removed from frame t . The flowchart of the object region tracking in backward direction is illustrated in Fig. 5, which is composed of the following five main steps.

Step 1: variable $bTemporal$ is used to indicate whether the corresponding foreground region of current foreground region continuously describes the same object in backward direction or not.

Step 2: each 4×4 block $B_t(k, l)$ in the foreground region is projected to the previous frame $t-1$ by using its MV (MV^x, MV^y) . Its projected position (k', l') in frame $t-1$ is

$$(k', l') = (k, l) + (MV^x, MV^y) \quad (4)$$

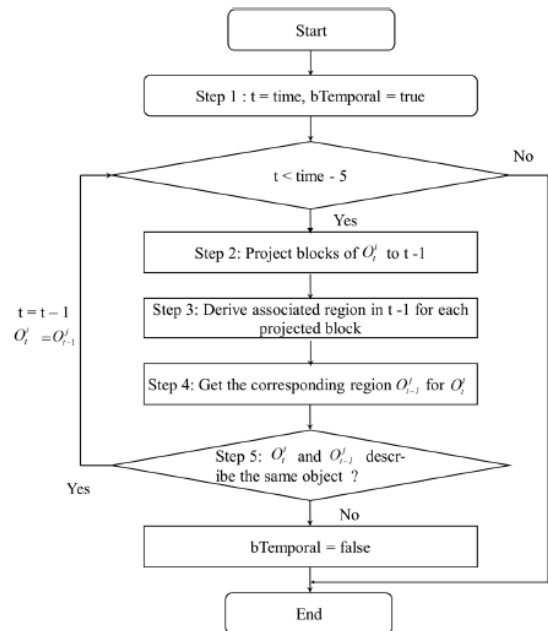


Fig. 5 The flowchart of object region tracking in backward direction.

Step 3: for each projected block $B_{t-1}(k', l')$ at time $t-1$, if $B_{t-1}(k', l')$ belongs to background region, the region for $B_{t-1}(k', l')$ will be set to ϕ .

Step 4: the number of projected blocks in each foreground region O_{t-1}^i at frame $t-1$ is counted. Then, the foreground region O_{t-1}^i with the most projected blocks will be identified as the corresponding foreground region for O_t^i at frame t . **Step 5:** the intersection of O_t^i and its corresponding foreground region O_{t-1}^i is computed and compared with an empirical threshold. In (5), $size(Object)$ denotes the number of foreground blocks in region $Object$ and is empirically set to 50%. If (5) is not satisfied, $bTemporal$ is set to be false and this process is finished. Otherwise, we assume and its corresponding foreground region is assumed to describe the same object region, and then return to step 1 and repeat the whole process.

$$\frac{size(O_t^i \cap O_{t-1}^i)}{size(O_t^i)} > \delta \quad (5)$$

B. Object Boundary Refinement

Figs. 6(a) and (b) show an example of block partitions for two surveillance video frames with moving persons and vehicles. Here, the largest square blocks with yellow border, smaller square blocks with green border and rectangular blocks with pink border represent coding tree unit (CTU), CU and PU respectively. It is observed that blocks within the moving person and vehicle region are encoded with smaller CU and PU sizes when compared to CU and PU sizes of the blocks within the background region. On the contrary, when the minimum size of CU and PU in one block row (column) is larger than the average size of CU and PU of the moving object region, this block row (column) has a high probability to belong to the background region. So we propose to refine the object boundary using the CU

and PU sizes. In order to describe the object boundary refinement more clearly, we use CU depth level in this section. The relationship between CU size and depth level is illustrated in Fig. 7. The depth range of CU is [0, 3]. The size of CTU in Fig. 7 is 64x64 and the root of CTU corresponds to depth level 0. The depth level for CU size of 32x32, 16x16 and 8x8 are 1, 2 and 3 respectively. One CU can be further split into one, two or four PUs according to the PU splitting type. When the CU is only split into one PU, the depth level for PU is 0; otherwise, the depth level of PU is 1. Note that all 4x4 blocks inside a CU share the same CU depth and PU depth. The depth of a 4x4 block is defined as the sum of CU and PU depth.

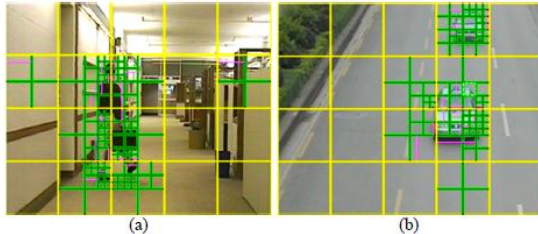


Fig. 6. Block partitions of the moving person and vehicles.

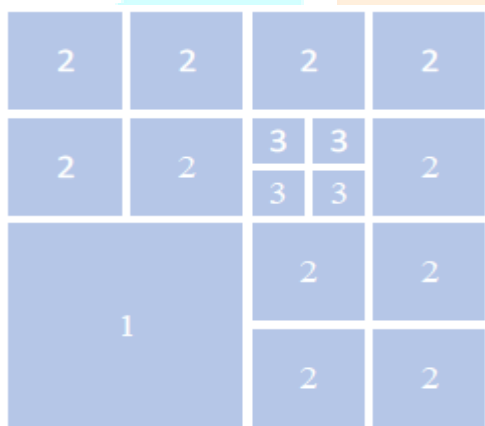


Fig. 7. The relationship between the depth level and CU size.

Our object boundary refinement is operated on the 4x4 block level, and the main steps of object boundary refinement is described as follows.

Step 1: the average depth of the foreground region is calculated. To be specific, let H_0 and H_1 denote the left boundary and right boundary of the foreground region, V_0 and V_1 denote the top and bottom boundary of the foreground region, $CUdepth(k, l)$ denotes the CU depth of the block at position (k, l) , $PUdepth(k, l)$ denotes the PU depth of the block at position (k, l) , $depth(k, l)$ denotes the depth of the block at position (k, l) and can be computed by (6), $foregBlockNum$ denotes the number of blocks within the foreground region, and $floor(x)$ represents the operation to get the smallest integer larger than or equal to x . Thus, the average depth of the foreground region is computed by (7).

$$depth(k, l) = CUdepth(k, l) + PUdepth(k, l) \tag{6}$$

$$AvgDepth = floor\left(\sum_{\substack{H_0 \leq k \leq H_1 \\ V_0 \leq l \leq V_1}} \frac{depth(k, l)}{foregBlockNum} + 0.5\right) - 1 \tag{7}$$

Step 2: the maximum depth for each block row within the moving object region is calculated and compared with $AvgDepth$. Formally, let $MaxDepth(k)$ denote the maximum depth in block row k , and it is computed by

$$MaxDepth(k) = \max_{H_0 \leq l \leq H_1} (depth(k, l)) \tag{8}$$

If (9) is satisfied, this block row will be marked as the candidate to be modified to background blocks.

$$MaxDepth(k) \leq AvgDepth \tag{9}$$

Step 3: for each candidate block row, we need to further check whether they are real background blocks or not. Formally, let $IndexMin(k)$ denote the index of the block with the minimum depth in row k , it is computed by

$$IndexMin(k) = \underset{l}{\operatorname{argmin}} (depth(k, l)) \tag{10}$$

The block with $IndexMin(k)$ is referred as the index block in row k . The CU which contains the index block is referred as the index CU. Then we get the top CU and the bottom CU of the index CU. After that, the maximum block row number in the top CU and the minimum block row number in the bottom CU are denoted as $MaxRowTop$ and $MinRowBottom$ respectively. If $MaxDepth(MaxRowTop)$ or $MaxDepth(MinRowBottom)$ is smaller than $AvgDepth$, the candidate row k has a high probability to belong to the background blocks. To be specific, if (11) or (12) is satisfied, we assume candidate row k belongs to background and need to be modified to background blocks.

$$MaxDepth(MaxRowTop) \leq AvgDepth \tag{11}$$

$$MaxDepth(MinRowBottom) \leq AvgDepth \tag{12}$$

Figs. 8 shows an example of the object boundary refinement on frame #27 of the Hall Monitor sequence. The $AvgDepth$ of the moving object region in Figs. 8(a) is 2. As illustrated in Figs. 8(b), the rows belonged to region 1 and region 2 satisfy (9), so the rows within region 1 and region 2 are marked as the candidates. In Figs. 8, the rows within region 1 neither satisfy (11) nor (12), so the rows within region 1 are foreground blocks. The rows within region 2 satisfy (12), so they are real background blocks and need to be removed from the foreground regions. After the boundary refinement, the moving object region are illustrated in Figs. 8(c). After checking each block row within the moving object, we flip the moving object region by 90 degree and use the same way to check each block column.

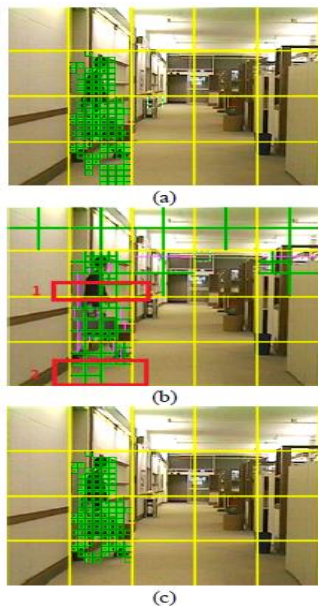


Fig. 8 (a) Moving object region without boundary refinement; (b) CU and PUsizes of the whole frame; (c) Moving object region after boundary refinement.

IV. MOVING OBJECT CLASSIFICATION IN HEVC COMPRESSED DOMAIN

For object classification in surveillance videos, we aim to classify the segmented moving objects into persons and vehicles using “bag of HEVC syntax words” in HEVC compressed domain. The “bag of words” representation has been successfully used for object classification in the pixel domain [26-27]. R. V. Badu et al. [28] propose to use “bag of words” representation in H.264/AVC compressed domain to classify the video content. The major contribution of this work is to establish a bag of words model in the HEVC domain for moving object classification. This proposed object classification has the following major

Steps:

- describing each coding block within the moving object region using HEVC syntax features;
- constructing a codebook using a clustering method;
- representing each moving object using a normalized histogram of codeword from the codebook; and
- training a binary classifier to classify the moving objects into persons and vehicles. The main challenge is to select effective features in the HEVC compressed domain.

In this work, we have identified three types of features, the length of motion vectors, prediction modes, and motion vector difference (MVD), as effective features for our object classification. The length of the motion vector relates to the velocity of the object, which is a simple yet important feature for discriminating persons and vehicles. This is because vehicles usually move faster than persons. It is observed that persons often undergo non-rigid deformations, it is harder to find a good match for each PU within the region of moving persons. As a result, more blocks within the region of persons are coded with intra modes when compared to the blocks within the region of moving vehicles. Therefore, the prediction mode can be utilized as an effective feature. For example, as shown in Figs. 9 (a) and (b), blocks

with red borders are encoded withintra mode, blocks with green borders are encoded with intermode and non-zero MVs, and other blocks are encoded with inter mode and zero MVs. We can see that more blocks are coded with intra modes within the moving person. Formally,

the prediction mode of current block at frame t is denoted as $CurrModet$. Because the MVs within vehicles are more consistent than those within persons, the MVD between neighboring blocks within the region of vehicles are usually smaller than that within the region of persons. Since we focus on the motion variation within the moving objects, the MVD of current block and its neighboring blocks is calculated only when they both have non-zero MVs. Specifically, MVD of current block is computed by

$$CurrMVD_t^i = \begin{cases} 0, & \text{if } NeighMV_t^i \text{ or } CurrMV_t = 0 \\ abs(NeighMV_t^i - CurrMV_t), & \text{else} \end{cases} \quad (13)$$

$$MaxCurrMVD_t = max(CurrMVD_t^i) \quad (14)$$

where $NeighMV_t^i$ denotes the MV of the i th neighboring block at frame t , $CurrMV_t$ denotes the MV of the current

block at frame t , $CurrMVD_t^i$ denotes the MVD between current block and the i th neighboring block, and $MaxCurrMVD_t$ denotes the maximum MVD between current block and its neighboring blocks at frame t . In addition, the maximum MVD of the collocated block at frame $t-1$ and $t+1$ are also used as the features of current block, which are denoted as $MaxCurrMVD_{t-1}$ and $MaxCurrMVD_{t+1}$ respectively.

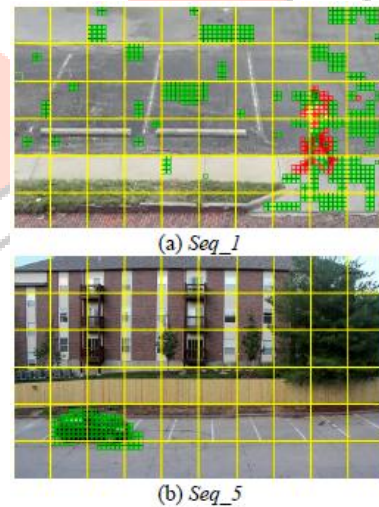


Fig. 9. The distribution of intra mode within the moving objects (blocks with red borders are encoded with intra mode, blocks with green borders are encoded with inter mode and non-zero MVs, and other blocks are encoded with intermode and zero MVs).

The features are summarized in Table I. The total feature size is 5. Once the features of all blocks in the training dataset have been extracted, these feature vectors are clustered into M clusters by using k -means clustering method. The center of each cluster becomes a codeword. In total, the codebook will have M codewords. For example, in our experiments, we set $M = 25$. Each moving object, either a person or a vehicle, will contain a large number of blocks. We compute the L_2 -norm distance between the feature vector of each block B and each codeword C_m . We find the

codeword which has the minimum distance to B_{nand} and cast B_{nto} to the bin of this codeword. In this way, we generate a codeword histogram for all blocks in the object. After normalized by its size, the histogram is used as the feature to describe the moving object. With this feature description scheme and the training data, we train a binary linear SVM classifier for person-vehicle classification. Our object classification algorithm is summarized in Algorithm 1.

TABLE I: Each Component Of The Feature Vector For Object Classification

Feature name	Definition	Feature Size
Motion vector	$CurrMV_t$	1
Prediction mode	$CurrMode_t$	1
MVD	$MaxCurrMVD_t$	1
	$MaxCurrMVD_{t-1}$	1
	$MaxCurrMVD_{t-2}$	1

Algorithm 1: Moving Object Classification

V. EXPERIMENTAL RESULTS

In order to train the person-vehicle model for moving object classification, 4 training sequences are used, which are illustrated in Figs. 10. To evaluate the performance of our proposed moving object segmentation and classification scheme in HEVC compressed domain, we have collected 2 sequences from CDNet2012 dataset (Highway and Pedestrians), 1 sequence from H.264/AVC standard sequence (Hall Monitor) and 6 sequences from our dataset.

Algorithm 1: Moving Object Classification

Input: Bounding box of the moving object, pre-trained codebook and classification model

Output: Classification result of the moving object, either person or vehicle

Begin

1. For each 4x4 image patch within the bounding box, form its feature vector according to Table I, and find the nearest codeword in the pre-trained codebook.
2. Generate a codeword histogram for all image patches within the bounding box.
3. Normalize the histogram by its size.
4. Use the pre-trained classification model to classify the moving object into person and vehicle.

End

4 of the test sequences have more than one objects in one frame, which are Highway, Seq_3, Seq_4, and Seq_6. In addition, there are persons and vehicles present in one frame in Seq_6. Example frames of the test videos are shown in Figs. 11 and Figs. 12. The resolutions and number of frames for the training and test videos are illustrated in Table II and Table III. Both the training and testing videos are encoded using the HEVC HM v10.0 encoder, at various bitrates, with the GOP structure IBBBB, i.e., the first frame is coded as intra (I), and subsequent frames are coded as generalized B frames. HEVC syntax features, such as motion vectors, prediction modes, CU sizes, and PU types, are extracted from HEVC compressed bitstream.

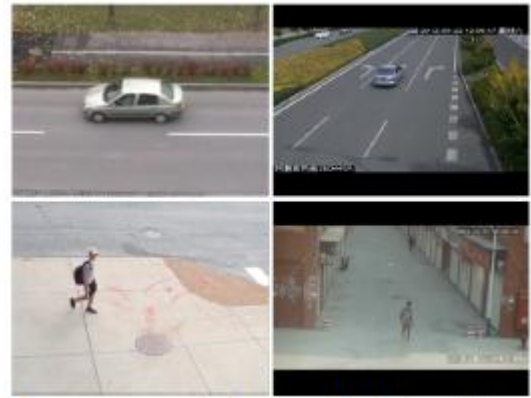


Fig. 10. Example frames of training videos.



Fig.11. Example frames of test videos from public dataset.

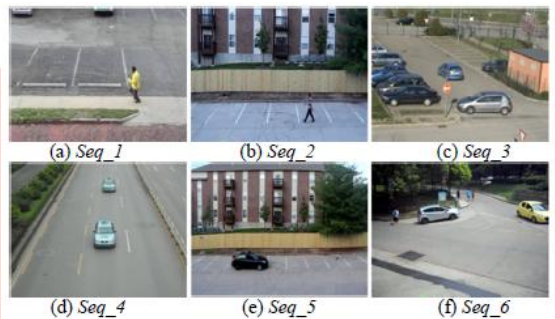


Fig. 12. Example frames of test videos from our dataset.

TABLE II: Resolutions And Number Of Frames For Each Training Sequence

Sequence	Resolution	Number of frames
Training_Seq_1	352x288	100
Training_Seq_2	640x480	200
Training_Seq_3	640x480	100
Training_Seq_4	640x480	100

TABLE III: Resolutions And Number Of Frames For Each Test Sequence

Sequence	Resolution	Number of frames
Hall Monitor	352x288	60
Highway	320x240	1100
Pedestrians	360x240	700
Seq_1	640x480	100
Seq_2	1920x1080	200
Seq_3	320x256	100
Seq_4	640x480	100
Seq_5	1920x1080	150
Seq_6	640x480	200

Fig. 13 shows several examples of the moving object segmentation results using our proposed method. As it is shown, the moving person or vehicle can be well detected and segmented when they are not close to each other whereas the moving person and vehicle will be segmented as one whole object when they are close to each other.

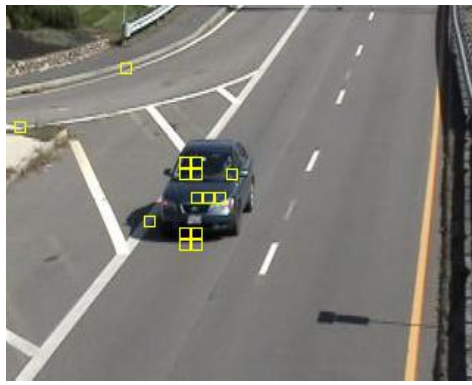


Fig. 13.a Preprocessing.



Fig. 13.e Output

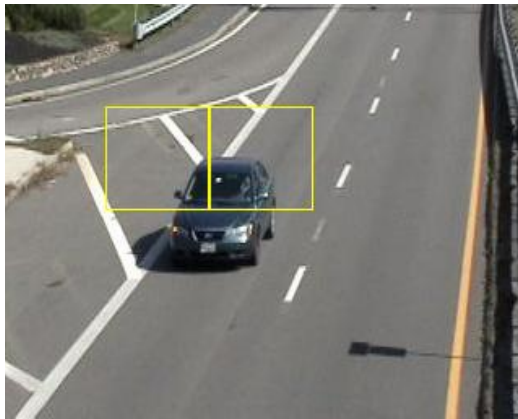


Fig. 13.b Object region tracking



Fig. 13.c Object Boundary Refinement



Fig. 13.d Object Segmentation

Moreover, the segmentation accuracy is measured by comparing the segmented foreground and background blocks with the ground truth labels for each frame of the test sequences. Specifically, the proposed moving object segmentation algorithm is evaluated in terms of precision, recall and Fmeasure, which are defined in (15) ~ (17). The notations TP, FP and FN are the total number of true positives, false positives, and false negatives respectively. Precision is defined as the number of TP divided by the total number of labeled 4x4 blocks. Recall is defined as the number of TP divided by the total number of ground truth labels. F-measure is the harmonic mean of precision and recall.

$$Precision = \frac{TP}{TP + FP} \quad (15)$$

$$Recall = \frac{TP}{TP + FN} \quad (16)$$

$$F - measure = \frac{2 \times Precision \times Recall}{Precision + Recall} \quad (17)$$

Table IV ~ Table VI summarize the performance of proposed segmentation algorithm in terms of Precision, Recall and Fmeasure for various values of QP between 22 and 32. It can be seen that the segmentation performance remains consistent for different QPs. In addition, for different sequences, these segmentation precision varies a lot. This is due to the different degree of motion vector noise produced from the encoding process for different sequences. As shown in Figs. 9, the motion vector noise of Seq_1 sequence is much larger than that of Seq_5 sequence. As a result, the segmentation performance of Seq_1 is relatively lower than Seq_5. Table VII compares the proposed method with the methods of [20], [14], and [10] in terms of precision, recall, and Fmeasure on Hall Monitor sequence. Note that the methods of [20], [14] and [10] are performed on H.264/AVC bitstream whereas our proposed method is performed on HEVC bitstream. This is because little research has been done on HEVC bitstream. When compared to these methods, our method can achieve similar segmentation performance as [20], which has the best performance among these three methods.

TABLE IV: Precision (In Percentage) For Various Qp Values

QP	Hall Monitor	Highway	Pedestrians	Seq_1	Seq_2	Seq_3	Seq_4	Seq_5	Seq_6	Avg
22	57.4	79.0	64.1	57.7	70.1	72.6	54.5	87.0	68.5	67.9
27	61.8	75.6	60.6	70.6	66.3	68.2	55.5	84.1	69.3	68.0
30	65.2	72.3	58.4	72.9	66.1	66.1	55.3	80.3	67.7	67.1
32	67.1	70.0	56.1	70.1	64.3	63.4	56.5	77.0	67.2	65.7

TABLE V: Recall (In Percentage) For Various Qp Values

QP	Hall Monitor	Highway	Pedestrians	Seq_1	Seq_2	Seq_3	Seq_4	Seq_5	Seq_6	Avg
22	92.2	90.2	80.0	94.1	92.0	97.0	98.4	93.0	93.5	92.3
27	88.8	89.8	83.2	90.8	92.1	96.1	96.6	93.3	92.7	91.5
30	86.9	89.8	81.4	88.8	90.2	95.5	94.6	94.1	89.6	90.1
32	85.4	90.8	79.9	88.3	89.1	92.9	93.4	93.5	90.1	89.3

TABLE VI: F-Measure (In Percentage) For Various Qp Values

QP	Hall Monitor	Highway	Pedestrians	Seq_1	Seq_2	Seq_3	Seq_4	Seq_5	Seq_6	Avg
22	70.8	84.2	71.2	71.5	79.6	83.0	70.1	89.9	79.1	77.7
27	72.9	82.1	70.1	79.4	77.1	79.8	70.5	88.5	79.3	77.7
30	74.5	80.1	68.0	80.1	76.3	78.1	69.8	86.7	77.1	76.7
32	75.2	79.1	65.9	78.2	74.7	75.4	70.4	84.5	77.0	75.6

At the meantime, the running speed of our segmentation method is much faster. For the test video with resolution 352x288, the processing speed of the method [20] is about 10 frames per second (fps), whereas the processing speed of our proposed method is over 400 fps. Note that our system is implemented by C language whereas the system of [20] is implemented by matlab. The processing speed of our proposed method for each sequence is illustrated in Table VIII. These results were obtained on a 2.30 GHz Intel Core i3 CPU with 8GB RAM.

TABLE VII: Comparison Of Several Methods In Terms Of Precision, Recall And F-Measure For Hall Monitor Sequence

Method	Precision	Recall	F-measure
Proposed	63.7	87.9	73.8
[20]	72.8	82.4	78.1
[14]	27.9	91.9	37.3
[10]	15.6	90.1	22.9

TABLE VIII: Running Speed Of The Proposed Moving Object Segmentation

Sequence	Resolution	Running Speed (fps)
Hall Monitor	352x288	437
Highway	320x240	645
Pedestrians	360x240	562
Seq_1	640x480	263
Seq_2	1920x1080	74
Seq_3	320x256	549
Seq_4	640x480	270
Seq_5	1920x1080	83
Seq_6	640x480	258

Table IX shows the comparison between the proposed segmentation method (preprocessing + moving object tracking and object boundary refinement) and the proposed method only with moving object tracking and object boundary refinement (without preprocessing) for all test sequences in terms of F-measure. As it is shown, the preprocessing process improves the F-measure accuracy about 14% for different QPs. In addition, Table X shows the comparison between the proposed segmentation method and the proposed segmentation only with preprocessing (without moving object tracking and object boundary refinement) for all test sequences in terms of F-measure. The moving object

tracking and object boundary refinement improves the F-measure accuracy about 5% for different QPs.

TABLE IX: Comparison Of The Proposed Method And The Proposed Method Without Preprocessing In Terms Of F-Measure

QP	Proposed segmentation	Proposed method without preprocessing process
22	77.7	61.2
27	77.7	63.1
30	76.7	63.8
32	75.6	64.6

TABLE X: Comparison Of The Proposed Segmentation And The Proposed Segmentation With Only Preprocessing In Terms Of F-Measure

QP	Proposed segmentation	Proposed segmentation with only preprocessing
22	77.7	73.5
27	77.7	72.6
30	76.7	71.9
32	75.6	70.5

TABLE XI: Accuracy Of The Proposed Classification Algorithm

QP	Hall Monitor	Highway	Pedestrians	Seq_1	Seq_2	Seq_3	Seq_4	Seq_5	Seq_6	Avg
22	94	71	88	99	93	85	84	89	89	88
27	96	76	94	99	94	91	94	95	86	92
30	94	80	83	96	85	94	94	87	84	89
32	91	81	87	94	86	95	93	95	88	90

Table XI shows the performance of person-vehicle classification results on different test videos and quantization settings. For performance evaluations, we manually label each moving object, either a person or a vehicle, as the ground truth. The moving object region containing just one kind of object is used for classification whereas the moving object region containing both the person and vehicle is not used for classification. In addition, the detected regions, which belong to background but falsely detected as moving object regions, are not used for classification. In total, there are 1187 persons and 1415 vehicles. We use the accuracy as the performance metric, which is defined as

$$Accuracy = \frac{TP}{TP + FP} \quad (18)$$

Here, TP and FP are true positive and false positive rates, respectively. We can see that the overall classification accuracy of the proposed method is about 90% for more than 2500 moving objects. Furthermore, our system achieves consistent performance across different QPs. Since the input of our classification algorithm is the moving object region extracted from the video frame by our segmentation algorithm, the classification accuracy of our classification algorithm depends on the segmentation accuracy of our segmentation algorithm. In Seq_3 and Seq_4, the objects are not close to each other, our classification algorithm can classify them with about 90% accuracy. In Highway sequence, the vehicles are near to each other, our segmentation algorithm fails to segment them correctly, which is illustrated in Figs. 13(c), this increases the difficulty for our classification algorithm. Even though, our algorithm can also classify them with about 80% accuracy. The average

running speed of the proposed segmentation and classification scheme for each sequence is illustrated in Table XII. Please note that the decoding time of MV and the associated mode information is not included. As it is shown, our proposed scheme can achieve more than 200fps for 640x480 sequence.

Table XIII compares our proposed method against pixel domain SVM classifier for person-vehicle classification using the same detected bounding boxes. In the pixel-domain SVM object classification, histograms of oriented gradients (HOG)[29] feature is used to describe each 4x4 patch within the moving object region. Then the moving object region is represented by "Bag of Words". We can see that our proposed method achieves comparable performance with the pixel domain SVM classifier. When compared to pixel domain SVM classification, our method directly uses the information from the bitstream and do not need the full decoding of the bitstream, which is a computation-intensive task. To verify the performance improvement of our proposed features in "bag of temporal-spatial words", we refer the system only including the motion vector of current block in "bag of temporal-spatial words" representation as the baseline system since the motion vector is the straightforward but discriminative feature to distinguish person and vehicles. Table XIII also compares our proposed system with the baseline system. It is revealed that the baseline system fails to distinguish the persons and vehicles when persons walk fast or vehicles move slowly, such as the fast-moving person in Seq_1 and Seq_2, and the slow-moving vehicles in Seq_5. This is because that most vehicles move faster than persons in the training dataset and the model trained from the dataset will misclassify the fast moving persons as vehicles. When compared to baseline system, our proposed system can distinguish these persons and vehicles robustly and efficiently.

TABLE XII: Processing Speed Of Proposed Segmentation And Classification Method

Sequence	Resolution	Running Speed (fps)
<i>Hall Monitor</i>	352x288	389
<i>Highway</i>	320x240	533
<i>Pedestrians</i>	360x240	457
<i>Seq_1</i>	640x480	241
<i>Seq_2</i>	1920x1080	62
<i>Seq_3</i>	320x256	435
<i>Seq_4</i>	640x480	233
<i>Seq_5</i>	1920x1080	71
<i>Seq_6</i>	640x480	233

TABLE XIII: Comparison Of Our Proposed Method With Pixel Domain SVM Classifier And The Baseline System For QP = 27

Sequence	Proposed	Pixel domain SVM classifier	Baseline system
<i>Hall Monitor</i>	96	97	91
<i>Highway</i>	76	93	61
<i>Pedestrians</i>	94	96	88
<i>Seq_1</i>	99	96	15
<i>Seq_2</i>	94	99	17
<i>Seq_3</i>	91	95	94
<i>Seq_4</i>	94	95	76
<i>Seq_5</i>	95	98	37
<i>Seq_6</i>	86	95	56
<i>Avg.</i>	92	96	60

VI. CONCLUSION

In this paper, we have presented a novel approach to segment and classify the moving objects from HEVC compressed surveillance video. Only the motion vectors and the associated coding modes from the compressed stream are used in the proposed method. In the proposed method, firstly, MV interpolation for intra-coded PU and MV outlier removal are employed for preprocessing. Secondly, blocks with non-zero motion vectors are clustered into connected foreground regions by the four-connectivity component labeling algorithm. Thirdly, object region tracking based on temporal consistency is applied to the connected foreground regions to remove the noise regions. The boundary of moving object region is further refined by the coding unit size and prediction unit size. Finally, a person-vehicle classification model using bags of spatial-temporal HEVC syntax words is trained to classify the moving objects, either persons or vehicles. The proposed method has a fairly low processing time, yet still provides high accuracy.

VII. REFERENCES

- [1] M. Grundmann, V. Kwatra, M. Han, and I. Essa, "Efficient hierarchical graph-based video segmentation," in Proc. IEEE Conf. Comput. Vis. and Pattern Recognit., pp. 2141–2148, Jun. 2010.
- [2] S. Chien, W. Chan, Y. Tseng, and H. Chen, "Video object segmentation and tracking framework with improved threshold decision and diffusion distance," IEEE Trans. Circuits Syst. Video Technol., vol. 23, no. 6, pp. 921–634, Jun. 2013.
- [3] H. Sakaino, "Video-based tracking, learning, and recognition method for multiple moving objects," IEEE Trans. Circuits Syst. Video Technol., vol. 14, no. 5, pp. 1661–1674, Oct. 2013.
- [4] "Generic Coding of Moving Pictures and Associated Audio Information- Part 2: Video," ITU-T and ISO/IEC JTC 1, ITU-T Recommendation H.262 and ISO/IEC 13 818-2 (MPEG-2), 1994.
- [5] T. Wiegand, G. J. Sullivan, G. Bjontegaard, and A. Luthra, "Overview of the H.264/AVC video coding standard," IEEE Trans. Circuits Syst. Video Technol., vol. 13, no. 7, pp. 560–576, Jul. 2003.
- [6] G. J. Sullivan, J.-R. Ohm, W.-J. Han, and T. Wiegand, "Overview of the High Efficiency Video Coding (HEVC) standard," IEEE Trans. Circuits Syst. Video Technol., vol. 22, no. 12, pp. 1649–1668, Dec. 2012.
- [7] R. V. Babu, M. Tom, and P. Wadekar, "A survey on compressed domain video analysis techniques," Multimedia Tools and Applications, vol. 75, no. 2, pp. 1043–1078, Jan. 2013.
- [8] R. V. Babu, K. R. Ramakrishnan, H. S. Srinivasan, "Video object segmentation: a compression domain approach," IEEE Trans. Circuits Syst. Video Technol., vol. 14, no. 4, pp. 462–474, Apr. 2004.
- [9] S. D. Bruyne, C. Poppe, S. Verstockt, P. Lambert, and R. V. D. Walle, "Estimating motion reliability to improve moving object detection in the H.264/AVC domain," in Proc. IEEE Int. Conf. Multimedia Expo., pp. 330–333, Jun. 2009.
- [10] W. Zeng, J. Du, W. Gao, and Q. Huang, "Robust moving object segmentation on H.264/AVC compressed video using the block-based MRF model," Real-Time Imaging, vol. 11, no. 4, pp. 290–299, Aug. 2005.

- [11] W. Lin, M. Sun, H. Li, Z. Chen, W. Li, and B. Zhou, "Macroblockclassification method for video applications involving motions," *IEEE Trans. Broadcasting*, vol. 58, no. 1, pp. 34–46, Mar. 2012.
- [12] Z. Liu, Y. Lu, and Z. Zhang, "Real-time spatiotemporal segmentation of video objects in the H.264 compression domain," *J. Visual Commun. Image Represent.*, vol. 18, no. 3, pp. 275–290, Jun. 2005.
- [13] Y. Chen and I. V. Bajic, "A joint approach to global motion estimation and motion segmentation from a coarsely sampled motion vector field," *IEEE Trans. Circuits Syst. Video Technol.*, vol. 21, no. 9, pp. 1316–1328, Sep. 2011.
- [14] Y. Chen, I. V. Bajic, and P. Saedi, "Moving region segmentation from compressed video using global motion estimation and Markov random fields," *IEEE Trans. Multimedia*, vol. 13, no. 3, pp. 421–431, Jun. 2011.
- [15] C. Poppe, S. D. Bruyne, T. Paridaens, P. Lambert, and R. V. D. Walle, "Moving object detection in the H.264/AVC compression domain for video surveillance applications," *J. Visual Commun. Image Represent.*, vol. 20, no. 6, pp. 428–437, Aug. 2009.
- [16] F. Porikli, F. Bashir, and H. Sun, "Compression domain Video Object Segmentation," *IEEE Trans. Circuits Syst. Video Technol.*, vol. 20, no. 1, pp. 2–14, Jan. 2010.
- [17] M. Laumer, P. Amon, A. Hutter, and A. Kaup, "Compressed domain moving object detection by spatio-temporal analysis of H.264/AVC syntax elements," in *Proc. IEEE Conf. Picture Coding Symposium (PCS)*, pp. 282–286, May 2015.
- [18] P. Dong, Y. Xia, L. Zhuo, and D. Feng, "Real-time moving object segmentation and tracking for H.264/AVC surveillance videos," in *Proc. IEEE Int. Conf. Image Processing, Brussels*, pp. 11–14, Sep. 2011.
- [19] H. Sabirin and M. Kim, "Moving object detection and tracking using a spatio-temporal graph in H.264/AVC bitstreams for video surveillance," *IEEE Trans. Multimedia*, vol. 14, no. 3, pp. 657–668, Jun. 2012.
- [20] S. H. Khatoonabadi and I. V. Bajic, "Video object tracking in the compression domain using spatio-temporal Markov random fields," *IEEE Trans. Image Process.*, vol. 22, no. 1, pp. 300–313, Jan. 2013.
- [21] J. Ohm, G. J. Sullivan, H. Schwarz, T. K. Tan, and T. Wiegand, "Comparison of the coding efficiency of video coding standards --including high efficiency video coding (HEVC)," *IEEE Trans. Circuits Syst. Video Technol.*, vol. 22, no. 12, pp. 1669–1684, Dec. 2012.
- [22] H. Li, Y. Zhang, M. Yang, Y. Men, and H. Chao, "A rapid abnormal event detection method for surveillance video based on a novel feature in compressed domain of HEVC," in *Proc. IEEE Int. Conf. Multimedia and Expo. (ICME)*, pp. 1–6, Jul. 2014.
- [23] B. Dey, and M. K. Kundu, "Efficient foreground extraction from HEVC compressed video for application to real-time analysis of surveillance 'big' data," *IEEE Trans. Image Process.*, vol. 24, no. 11, pp. 3574–3585, Nov. 2015.
- [24] D. Park, D. Lee, and S. Oh, "Object tracking in HEVC bitstreams," *Journal of Broadcast Engineering*, vol. 20, no. 3, pp. 449–463, May 2015.
- [25] L. He, Y. Chao, and K. Suzuki, "A run-based two-scan labeling algorithm," *IEEE Trans. Image Process.*, vol. 17, no. 5, pp. 749–756, May 2008.
- [26] G. Csurka, C. Dance, L. Fan, J. Willamowski, and C. Bray. "Visual categorization with bags of keypoints," In *Workshop on statistical learning in computer vision, ECCV*, pp. 1–22, May 2004.
- [27] F. Wang, Z. Sun, Y. Jiang, and C. Ngo, "Video event detection using motion relativity and feature selection," *IEEE Trans. Multimedia*, vol. 16, no. 5, pp. 1303–1315, Aug. 2014.
- [28] S. Biswas, and R. V. Babu. "H. 264 compressed video classification using histogram of oriented motion vectors (HOMV) ," in *Proc. IEEE Conf. Acoustics, Speech and Signal Processing (ICASSP)*, pp. 2040–2044, May 2013.
- [29] N. Dalal, and B. Triggs, "Histograms of oriented gradients for human detection," in *Proc. IEEE Conf. Comput. Vis. and Pattern Recognit.*, pp. 886–893, Jun. 2005.
- [30] M. Moriyama, K. Minemura and K.S. Wong, "Moving object detection in HEVC video by frame sub-sampling," *IEEE Int. Symposium on ISPACS*, pp. 48–53, Nov. 2015.
- [31] S. Gil, J. T. Meyer, T. Schierl, C. Hellge, and W. Samek, "Hybrid video object tracking in H.265/HEVC video streams," *IEEE Int. conf on MMSP*, 2016.
- [32] S. Pulare and S. S. Tale, "Implementation of Real Time Multiple Object Detection and Classification of HEVC Videos," *Int. Journal for IRST*, vol. 2, no. 11, pp. 248–254, 2016.

