# Discriminative Video Pattern Search for Efficient Action Detection

Junsong Yuan, Member, IEEE, Zicheng Liu, Senior Member, IEEE, Ying Wu, Senior Member, IEEE

Abstract-Actions are spatio-temporal patterns. Similar to the sliding window-based object detection, action detection finds the reoccurrences of such spatio-temporal patterns through pattern matching, by handling cluttered and dynamic backgrounds and other types of action variations. We address two critical issues in pattern matching-based action detection: (1) the intra-pattern variations in actions, and (2) the computational efficiency in performing action pattern search in cluttered scenes. First, we propose a discriminative pattern matching criterion for action classification, called naive-Bayes mutual information maximization (NBMIM). Each action is characterized by a collection of spatio-temporal invariant features and we match it with an action class by measuring the mutual information between them. Based on this matching criterion, action detection is to localize a subvolume in the volumetric video space that has the maximum mutual information toward a specific action class. A novel spatio-temporal branch-and-bound (STBB) search algorithm is designed to efficiently find the optimal solution. Our proposed action detection method does not rely on the results of human detection, tracking or background subtraction. It can well handle action variations such as performing speed and style variations, as well as scale changes. It is also insensitive to dynamic and cluttered backgrounds and even to partial occlusions. The crossdataset experiments on action detection, including KTH, CMU action datasets, and another new MSR action dataset, demonstrate the effectiveness and efficiency of the proposed multi-class multipleinstance action detection method.

Index Terms—video pattern search, action detection, spatiotemporal branch-and-bound search

# **1** INTRODUCTION

Detecting human actions in video sequences is an interesting yet challenging problem. It has a wide range of applications including video surveillance, tele-monitoring of patients and senior people, medical diagnosis and training, video indexing, and intelligent human computer interaction, etc. Actions can be treated as spatio-temporal objects which are characterized as spatio-temporal volumetric data. Like the use of sliding windows in object detection, action detection can be formulated as locating spatio-temporal subvolumes in videos (*i.e.* video patterns) that contain the target actions. Despite previous successes of sliding window-based object detection [1] [2], this approach cannot easily be extended to action detection. It is still a challenging problem to detect and locate actions in video sequences, mainly due to the following two difficulties.

First, the computational complexity of pattern searching in the video space is much higher than that of object search in the image space. Without any prior knowledge about the location, temporal duration, and the spatial scale of the action, the search space for video patterns is prohibitive for exhaustive search. For example, a oneminute video sequence of size  $160 \times 120 \times 1800$  contains billions of valid spatio-temporal subvolumes of various sizes and locations. Therefore, although the state-of-the-art approaches of object detection can efficiently search the spatial image space [3] [1], they are in general not scalable to search videos, due to such an enormous search space. To reduce this huge search space, some other methods try to avoid exhaustive search by sampling the search space, e.g. only considering a fixed number of spatial and temporal scales [4]. However, this treatment is likely to cause missing detections. Moreover, the solution space is still quite large even after subsampling.

Second, human actions often exhibit tremendous amount of intra-pattern variations. The same type of actions may look very different in their visual appearances. There are many factors that contribute to such variations including the performing speed, clothing, scale, view points, not to mention partial occlusions and cluttered backgrounds. When using a single and rigid action template for pattern matching as in [4] [5], the actions that vary from the template cannot be detected. One potential remedy is to use multiple templates to cover more variations, but the required number of templates will increase rapidly, resulting in formidable computational costs.

We propose an efficient action detection approach that addresses these two challenges mentioned above. Each action is characterized by a set of spatio-temporal interest points (STIP) [6]. Provided a test video sequence, each STIP casts a positive or negative-valued vote for the action class, based on its mutual information with respect to that action class. As illustrated in Fig. 1, detection of an action is to search for a spatio-temporal subvolume that has the maximum total vote. Such a subvolume with maximum voting score also maximizes the point-wise mutual infor-

Junsong Yuan is with the School of Electrical and Electronics Engineering, Nanyang Technological University, 639798 Singapore (e-mail: jsyuan@ntu.edu.sg).

Zicheng Liu is with Microsoft Research, Redmond, WA, 98052 USA (email: zliu@microsoft.com.)

Ying Wu is with the department of Electrical Engineering and Computer Science, Northwestern University, Evanston, IL, 60208 USA (e-mail: yingwu@ece.northwestern.edu).



Fig. 1. Action detection through spatio-temporal pattern search. The highlighted subvolume has the maximum mutual information toward the specific event class. Each circle represents a spatio-temporal feature point which contributes a vote based on its own mutual information. The whole video volume is of size  $m \times n \times t$ , where  $m \times n$  is the image size and *t* is the temporal length.

mation toward a specific action class. Thus it is treated as a valid detection instance of that action class. This is a new formulation of action detection.

To handle the intra-class action variations, various action templates that belong to the same class (or pattern) are collected and the collection of all the STIPs forms a pool of positive STIPs. As each action pattern is represented by more than one template, it is able to tolerate the variations in actions. In terms of pattern matching, to make an analogy to the template-based pattern matching where only positive templates are utilized, our pattern matching is called *discriminative matching* because of the use of both positive and negative templates. Given both the positive and negative training STIPs, a classification scheme called naive-Bayes mutual information maximization (NBMIM) is proposed to classify a query video clip, i.e., a cloud of STIPs. By exploring discriminative learning, such a discriminative matching method can better distinguish the target action patterns from the cluttered and the moving backgrounds. Thus a more robust pattern matching can be achieved.

To handle the large search space in video, we propose a method that decouples the temporal and spatial spaces and applies different search strategies, respectively. By combining dynamic programming in the temporal space and the branch-and-bound in the spatial space, the proposed spatiotemporal branch-and-bound (STBB) method significantly speeds up the search of the spatio-temporal action patterns. Moreover, we also investigate how to detect multiple action instances simultaneously, *i.e.*, search for multiple subvolumes whose scores are higher than the detection threshold. Based on the new scheme, it can terminate many unnecessary candidates earlier during the process of branch-and-bound to save computation. It leads to a much faster search, without significantly degrading the quality of the detection results.

The benefits of our new method are three-fold. First, the proposed discriminative pattern matching can well handle action variations by leveraging all of the training data instead of a single template. By incorporating the negative training information, our pattern matching has stronger discriminative power across different action classes. Second, our method does not rely on object tracking, detection, and background subtraction. It can handle background clutters and other moving objects in the background. Last but not the least, the proposed spatio-temporal branch-and2

action detection method, it has been tested on various data sets, including cross-dataset experiments where the positive training data, negative training data, and test data are from different sources. The action categorization results on the KTH dataset are comparable to the state-of-the-art results. The multi-class multiple-instance action detection results demonstrate the effectiveness and efficiency of our method.

# 2 RELATED WORK

# 2.1 Action Recognition

Action recognition has been an active research topic. Given a video sequence, it requires to identify which type of action is performed in this video. Some previous works perform human action recognition based on the tracking of human body parts. For example, motion trajectories are used to represent actions in [7] [8] [9] [10] [11]. Unfortunately, robust object tracking is itself a non-trivial task. The problem becomes particularly difficult when there are occlusions or when the background is cluttered. Instead of relying on the body part information, some approaches use the silhouette information, such as using key poses in representing actions [12] [13]. Some other approaches treat actions as spatio-temporal shapes and characterize them by using manifold invariants in the spatio-temporal space. For example, in [14] a spatio-temporal volume (STV) is generated by the 2-D contours of the object along the temporal axis. By considering the STV as a 3-D manifold, this method extracts algebraic invariants of the manifold, which correspond to the changes in direction, speed and shape of the parts. Space-time shapes are also applied in [15]. It utilizes the properties of the solution to the Poisson equation to extract space-time features. Because both silhouettes and spatio-temporal shapes can be obtained through foregroundbackground separation, such action recognition approaches perform well when the background is reasonably clean or static. If the background is cluttered and dynamic, extracting foreground becomes very difficult. The noisy and erroneous silhouettes or spatio-temporal shapes largely limit the performance of these methods.

To avoid the foreground-background separation, many recent methods applied local spatio-temporal features to characterize actions and perform action classification over the set of local features [16] [17] [18] [19] [20] [21] [22] [23] [24]. In [19], spatio-temporal interest points (STIP) are proposed and applied to characterize human actions. In [18], local spatio-temporal features are quantized into "visual words" and the support vector machine is applied for classification. In [25], a video sequence is characterized by a "bag-of-words", where each frame corresponds to a "word". A semi-latent topic model is then trained for action recognition. These previous good recognition results validate the advantages of using the spatio-temporal local features.

Besides using local features, there are many previous works in designing and fusing various types of features for classification. In [26], a set of kinematic features are proposed for action recognition. In [27], mid-level motion features are developed from low-level optical flow information for action recognition. In [28], a motion descriptor based on optical flow measurements in a spatio-temporal volume is introduced. All of these methods require optical flow estimation. In [29], a motion-based representation called motion context is proposed for action recognition. To improve the recognition performance, both shape and motion information are used for action detection [30] [31] [32]. In [33], multiple features are fused using Fiedler Embedding. To select informative features, PageRank algorithm is applied in [17]. Based on the Gaussian processes with multiple kernel covariance functions, [34] proposes a Bayesian classification method to automatically select and weight multiple features. Spatio-temporal context information is utilized in [11] to improve the performance of action recognition. In [35], a generative model is learned by using both semantic and structure information for action recognition and detection.

#### 2.2 Action Detection

Different from action recognition or categorization [36] [19] [37], where each action video is classified into one of the pre-defined action classes, the task of action detection [38] [39] [4] [40] [41] needs to identify not only which type of action occurs, but also where (spatial location in the image) and when (temporal location) it occurs in the video. As discussed in [42], it is in general a more challenging problem as it needs to not only recognize the action, but also to locate it in the video space [39] [4] [43] [40] [41] [44]. Some previous methods apply template-based action matching [5] [39] [45]. For example, two types of temporal templates are proposed in [5] for characterizing actions: (1) the motion energy image (MEI); and (2) the motion history image (MHI). To provide a viewpoint-free representation for human actions, [46] introduces motion history volume (MHV). Besides motion templates, some other approaches also characterize an action as a sequence of postures, so that sequence matching methods can be applied to action recognition and detection [47] [48] [32]. In general, template-based approach is sensitive to the cluttered and dynamic backgrounds. To address this problem, [4] proposes to over-segment the video into many spatio-temporal video volumes. An action template is then matched by searching among these over-segmented video volumes. However, because only one template is utilized in matching, previous templates-based methods also have difficulties in handling intra-class action variations.

In addition to template matching, discriminative learning methods have also been applied to action detection. Motivated by the successful face detector [2], [38] extends the Haar features to the 3-dimensional space, followed by the boosting algorithm to integrate these features for classification. In [41], multiple instance learning is presented for learning a human action detector. However, head detection and tracking are required to help locate the person. To learn the action representations, [49] proposes to collect images from the Web and use this knowledge to automatically annotate actions in videos. In [50], boosted space-time window classifiers are introduced to detect human actions on real movies with substantial variation of actions in terms of subject. Both human motion and shape features are applied. As an earlier version of this paper, [40] proposes an efficient 3-dimensional branch-and-bound search for efficient action detection. This method is further developed in [51] for transductive action detection, where the requirement of training labels is reduced.

#### 2.3 Object Recognition and Detection

Besides action recognition and detection, some recent works in object recognition and detection were also related to our work. A recent work in [52] proposed the naive-Bayes nearest-neighbor (NBNN) classifier for image classification. In [1], object detection is formulated as finding the optimal bounding box that gives the highest detection score in the image. An efficient branch-and-bound method is proposed to search for the optimal bounding box in the image. Despite the successful applications in object detection [1] and image retrieval [53], it still is a nontrivial problem to extend the efficient search method from the spatial image space to the spatio-temporal video space. Thus, a further study is required.

# **3** CLASSIFICATION MODEL OF ACTIONS

#### 3.1 Interest Point Representation for Actions

We represent an action as a space-time object and characterize it by a collection of spatio-temporal interest points (STIPs) [6]. Two types of features are used to describe the STIPs [19]: histogram of gradient (HOG) and histogram of flow (HOF), where HOG is the appearance feature and HOF is the motion feature. These features have showed promising results in action categorization [19]. We denote a video sequence by  $\mathcal{V} = {\mathbf{I}_t}$ , where each frame  $\mathbf{I}_t$  consists of a collection of STIPs. We do not select key-frames but collect all STIPs to represent a video clip by  $\mathcal{Q} = {d_i}$ .

# 3.2 Naive-Bayes Mutual Information Maximization

We denote by  $d \in \mathbb{R}^N$  a feature vector describing a STIP and by  $\mathbf{C} \in \{1, 2, ..., C\}$  a class label. Based on the naive-Bayes assumption and by assuming the independence among the STIPs, we can evaluate the *pointwise mutual information* between a video clip  $\mathcal{Q}$  and a specific class  $c \in \{1, 2, ..., C\}$  as:

$$MI(\mathbf{C} = c, \mathcal{Q})$$

$$= \log \frac{P(\mathcal{Q}|\mathbf{C} = c)}{P(\mathcal{Q})} = \log \frac{\prod_{d_q \in \mathcal{Q}} P(d_q|\mathbf{C} = c)}{\prod_{d_q \in \mathcal{Q}} P(d_q)}$$

$$= \sum_{d_q \in \mathcal{Q}} \log \frac{P(d_q|\mathbf{C} = c)}{P(d_q)} = \sum_{d_q \in \mathcal{Q}} s^c(d_q), \quad (1)$$

where  $s^{c}(d_{q}) = MI(\mathbf{C} = c, d_{q})$  is the pointwise mutual information score to measure the association between  $d_{q}$  and class c. Assuming the independence among  $d_q$ , the final decision of Q is based on the summation of the mutual information from all primitive features  $d_q \in Q$  w.r.t. class c.

To evaluate the contribution  $s^c(d_q)$  of each  $d_q \in Q$ , we develop the pointwise mutual information through discriminative learning [54]:

$$s^{c}(d_{q}) = MI(\mathbf{C} = c, d_{q}) = \log \frac{P(d_{q}|\mathbf{C} = c)}{P(d_{q})}$$

$$= \log \frac{P(d_{q}|\mathbf{C} = c)}{P(d_{q}|\mathbf{C} = c)P(\mathbf{C} = c) + P(d_{q}|\mathbf{C} \neq c)P(\mathbf{C} \neq c)}$$

$$= \log \frac{1}{P(\mathbf{C} = c) + \frac{P(d_{q}|\mathbf{C} \neq c)}{P(d_{q}|\mathbf{C} = c)}P(\mathbf{C} \neq c)}.$$
(2)

If the prior probabilities are equal, *i.e.*  $P(\mathbf{C} = c) = \frac{1}{C}$ , we further have:

$$s^{c}(d_{q}) = \log \frac{C}{1 + \frac{P(d_{q}|\mathbf{C}\neq c)}{P(d_{q}|\mathbf{C}=c)}(C-1)}.$$
(3)

From Eq. 3, we can see that the likelihood ratio test  $\frac{P(d_q|\mathbf{C}\neq c)}{P(d_q|\mathbf{C}=c)}$  determines whether  $d_q$  votes positively or negatively for the class c. When  $MI(\mathbf{C} = c, d_q) > 0$ , *i.e.* likelihood ratio  $\frac{P(d_q|\mathbf{C}\neq c)}{P(d_q|\mathbf{C}=c)} < 1$ ,  $d_q$  votes a positive score  $s^c(d_q)$  for the class c. Otherwise if  $MI(\mathbf{C} = c, d_q) \leq 0$ , *i.e.*  $\frac{P(d_q|\mathbf{C}\neq c)}{P(d_q|\mathbf{C}=c)} \geq 1$ ,  $d_q$  votes a negative score for the class c. After receiving the votes from every  $d_q \in \mathcal{Q}$ , we can make the final classification decision for  $\mathcal{Q}$  based on its mutual information toward C classes.

For the C-class action categorization, we build C oneagainst-all classifiers. The test action Q is classified as the class that gives the maximum detection score.

$$c^* = \arg \max_{c \in \{1,2,..,C\}} MI(c, \mathcal{Q}) = \arg \max_{c \in \{1,2,..,C\}} \sum_{d \in \mathcal{Q}} s^c(d).$$

We call this naive-Bayes mutual information maximization (NBMIM). Compared with the naive-Bayes nearestneighbor (NBNN) [52], each score  $s^c(d)$  corresponds to the pointwise mutual information and can either be positive or negative. As will be explained in Section 4, such a property brings extra benefits in formulating action detection as a subvolume search problem, where a computationally efficient detection solution can be found.

#### 3.3 Likelihood Ratio Measurement

Denote by  $\mathbb{T}^{c+} = \{\mathcal{V}_i\}$  the positive training dataset of class c, where  $\mathcal{V}_i \in \mathbb{T}^{c+}$  is a video of class c. As each  $\mathcal{V}$  is characterized by a collection of STIPs, we represent the positive training data by the collection of all positive STIPs:  $\mathbb{T}^{c+} = \{d_j\}$ . Symmetrically, the negative data is denoted by  $\mathbb{T}^{c-}$ , which is the collection of all negative STIPs.

To evaluate the likelihood ratio for each  $d \in Q$ , we apply the Gaussian kernel density estimation based on the training data  $\mathbb{T}^{c+}$  and  $\mathbb{T}^{c-}$ . With a Gaussian kernel

$$K(d - d_j) = \frac{1}{\sqrt{2\sigma}} \exp^{-\frac{1}{2\sigma^2} \|d - d_j\|^2}$$

we adopt the nearest neighbor approximation as in [52]. The likelihood ratio becomes:

$$\frac{P(d|\mathbf{C} \neq c)}{P(d|\mathbf{C} = c)} = \frac{\frac{1}{|\mathbb{T}^{c-1}|} \sum_{d_j \in \mathbb{T}^{c-1}} K(d - d_j)}{\frac{1}{|\mathbb{T}^{c+1}|} \sum_{d_j \in \mathbb{T}^{c+1}} K(d - d_j)} \approx \exp^{-\frac{1}{2\sigma^2} (\|d - d_{NN}^{c--}\|^2 - \|d - d_{NN}^{c++}\|^2)} .$$
(4)

Here  $d_{NN}^{c-}$  and  $d_{NN}^{c+}$  are the nearest neighbors of din  $\mathbb{T}^{c+}$  and  $\mathbb{T}^{c-}$ , respectively. We approximate the numerator  $\frac{1}{|\mathbb{T}^{c-}|} \sum_{d_j \in \mathbb{T}^{c-}} K(d-d_j)$  by  $\exp^{-\frac{1}{2\sigma^2} ||d-d_{NN}^{c-}||^2}$ , and the denominator  $\frac{1}{|\mathbb{T}^{c+}|} \sum_{d_j \in \mathbb{T}^{c+}} K(d-d_j)$  by

In kernel based density estimation, it is difficult to select an appropriate kernel bandwidth  $\sigma$ . A large kernel bandwidth may over-smooth the density estimation, while a too small kernel bandwidth only counts on the nearest neighbors in the Parzen estimator. Let

$$\gamma(d) = \|d - d_{NN}^{c-}\|^2 - \|d - d_{NN}^{c+}\|^2.$$
(5)

According to Eq. 3 and Eq. 4, a positive  $\gamma(d)$  will generate a positive score  $s^c(d_q)$ , while a negative  $\gamma(d)$  will generate a negative  $s^c(d)$ . To avoid the selection of the best bandwidth  $\sigma$ , we adaptively adjust  $\sigma$  based on the *purity* in the neighborhood of a STIP d:

$$\frac{1}{2\sigma^2} = \begin{cases} \frac{|NN_{\epsilon}^{\epsilon+}(d)|}{|NN_{\epsilon}(d)|}, & if \quad \gamma(d) \ge 0\\ \frac{|NN_{\epsilon}^{\epsilon-}(d)|}{|NN_{\epsilon}(d)|}, & if \quad \gamma(d) < 0 \end{cases},$$
(6)

where  $NN_{\epsilon}^{c+}(d) = \{d_j \in \mathbb{T}^{c+} : \|d_j - d\| \leq \epsilon\}$  is the  $\epsilon$ -nearest neighbors of point d in the positive class  $c; NN_{\epsilon}^{c-}(d) = \{d_j \in \mathbb{T}^{c-} : \|d_j - d\| \leq \epsilon\}$  is the  $\epsilon$ -nearest neighbors of point d in the negative class;  $NN_{\epsilon}(d) = \{d_j \in \mathbb{T}^{c+} \cup \mathbb{T}^{c-} : \|d_j - d\| \leq \epsilon\}$  is the entire set of  $\epsilon$ -nearest neighbors of d. With  $\gamma(d)$  determining the sign of the vote  $s^c(d)$ , if d is located in a high purity region of the corresponding class, its vote  $s^c(d)$  is stronger.

#### 3.3.1 Efficient Nearest Neighbor Search

To obtain the voting score  $s^c(d)$ , for each  $d \in Q$ , we need to search for its nearest neighbors (NNs). To improve the efficiency of searching in the high-dimensional feature space and to obtain  $NN^c_{\epsilon}(d)$  quickly, we employ locality sensitive hashing (LSH) [55] to perform the approximate  $\epsilon$ -nearest neighbors ( $\epsilon$ -NN) search.

Based on  $NN_{\epsilon}^{c+}(d)$  and  $NN_{\epsilon}^{c-}(d)$ , instead of searching for the global nearest neighbor for each class, we approximate it by the closest point to the query d in the  $\epsilon$ -NN set. Taking the negative class as an example, we have:

$$||d - d_{NN}^{c-}|| = \min_{d_i \in NN_{\epsilon}^{c-}(d)} ||d - d_i||.$$

It is worth noting that  $d_{NN}^{c-}$  depends on the selection of  $\epsilon$ . If we happen to have  $|NN_{\epsilon}^{c-}(d)| = 0$ , we assume the negative nearest neighbor is at distance  $\epsilon$ , namely  $||d - d_{NN}^{c-}||^2 = \epsilon^2$  in Eq. 5. Applying the same strategy to the positive class, we have:

$$|d - d_{NN}^{c+}|| = \min_{d_i \in NN_{\epsilon}^{c+}(d)} ||d - d_i||.$$

When  $|NN_{\epsilon}^{c+}(d)| = 0$ , we assume  $||d - d_{NN}^{c+}||^2 = \epsilon^2$  in Eq. 5.

# 4 DISCRIMINATIVE VIDEO PATTERN SEARCH

Based on the proposed NBMIM criterion, action detection is to find a subvolume of maximum mutual information. As illustrated in Fig. 1, given a video sequence  $\mathcal{V}$ , we want to find a spatio-temporal subvolume  $V^* \subset \mathcal{V}$  with the highest mutual information score. Since STIPs are sparse features and involve only a very small number of pixels  $d \in \mathcal{V}$ , the optimal subvolume  $V^*$  may not be a unique one. For example, if a frame does not contain any STIPs, it becomes arbitrary for  $V^*$  to include this empty frame, as it does not affect the total voting score. To avoid this problem, we introduce a very small negative vote  $s(d_{\emptyset}) < 0$  to the empty pixels that are not associated with any STIP. Such a negative prior discourages the inclusion of empty pixels into  $V^*$ .

Given a specific class *c*, our target is to search for the optimal subvolume:

$$V^* = \arg \max_{V \subseteq \mathcal{V}} MI(V, \mathbf{C} = c)$$
(7)  
= 
$$\arg \max_{V \subseteq \mathcal{V}} \sum_{d \in V} s^c(d) = \arg \max_{V \in \mathbf{\Lambda}} f(V),$$

where  $f(V) = \sum_{d \in V} s^c(d)$  is the objective function and  $\Lambda$  denotes the candidate set of all valid subvolumes in  $\mathcal{V}$ . Suppose the target video  $\mathcal{V}$  is of size  $m \times n \times t$ . The optimal solution  $V^* = t^* \times b^* \times l^* \times r^* \times s^* \times e^*$  has 6 parameters to be determined, where  $t^*, b^* \in [0, m]$  denote the top and bottom positions,  $l^*, r^* \in [0, n]$  denote the left and right positions, and  $s^*, e^* \in [0, t]$  denote the start and end positions. As a counterpart of the bounding-box based object detection, the solution  $V^*$  is the 3D bounding volume that has the highest score for the target action.

The total number of the subvolumes is in the order of  $O(n^2m^2t^2)$ . Therefore, it is computationally prohibitive to perform an exhaustive search to find the optimal subvolume  $V^*$  from such an enormous candidate pool. In the following, we first present the conventional branch-and-bound solution extended directly from 2D bounding-box search in [1], and then present our new method to find  $V^*$  more efficiently.

#### 4.1 Spatio-Temporal Branch-and-Bound Search

# 4.1.1 Conventional branch-and-bound search

A branch-and-bound solution is proposed in [1] for searching the optimal bounding box in an image for object detection. This idea can be directly extended to find the optimal subvolume in videos, by replacing the spatial bounding box by a spatio-temporal subvolume.

Denote by  $\mathbb{V}$  a collection of subvolumes. Assume there exist two subvolumes  $V_{min}$  and  $V_{max}$  such that for any  $V \in \mathbb{V}$ ,  $V_{min} \subseteq V \subseteq V_{max}$ . Then we have

$$f(V) \le f^+(V_{max}) + f^-(V_{min}),$$
 (8)

where  $f^+(V) = \sum_{d \in V} \max(s^c(d), 0)$  contains the positive votes, while  $f^-(V) = \sum_{d \in V} \min(s^c(d), 0)$  contains the

negative ones. We denote the upper bound of f(V) for all  $V \in \mathbb{V}$  by:

$$\hat{f}(\mathbb{V}) = f^+(V_{max}) + f^-(V_{min}) \ge \max_{V \in \mathbb{V}} f(V).$$
(9)

Moreover, it is easy to see that if V is the only element in  $\mathbb{V}$ , we have the equality:

$$\hat{f}(\mathbb{V}) = f(V). \tag{10}$$

Eq. 9 and Eq. 10 thus meet the two requirements discussed in [1] for the effective upper bound in the branch-andbound search. With the first condition in Eq. 9,  $\hat{f}(\mathbb{V})$  is an upper bound of f(V). Therefore, it does not incur miss detection by using  $\hat{f}(\mathbb{V})$  for pruning unsatisfactory candidates. It guarantees the optimality of the solution. The second condition in Eq. 10 provides the termination condition of the branch-and-bound.

1	Algorithm 1: Conventional branch-and-bound (BB)
5	search (extension of [1])
	input : video $\mathcal{V} \in \mathbb{R}^{n \times m \times t}$ ; quality bounding function $\hat{f}$ (see text) output : $V^* = \arg \max_{V \subseteq \mathcal{V}} f(V)$
1	initialize $P$ as an empty priority queue
2	set $\mathbb{V} = [0, n] \times [0, n] \times [0, m] \times [0, m] \times [0, t] \times [0, t]$
3	while $\mathbb{V}$ contains more than one element <b>do</b>
4	split $\mathbb{V} \to \mathbb{V}^1 \cup \mathbb{V}^2$
5	get upper bound $\hat{f}(\mathbb{V}^1)$
6	push ( $\mathbb{V}^1$ , $\hat{f}(\mathbb{V}^1)$ ) into P
7	get upper bound $\hat{f}(\mathbb{V}^2)$
8	push ( $\mathbb{V}^2$ , $\hat{f}(\mathbb{V}^2)$ ) into $P$
9	retrieve top state $\mathbb{V}$ from $P$ based on $\widehat{f}(\mathbb{V})$
10	return $V^* = \mathbb{V}$

In order to distinguish this method from our new method, we call it conventional branch-and-bound method (Alg. 1). Compared to the spatial bounding box searching, the search of spatio-temporal subvolume is much more difficult. In videos, the search space has two additional parameters (the start and end on the time dimension) and expands from 4 dimensions to 6 dimensions. As the complexity of the branch-and-bound grows exponentially with respect to the number of dimensions, the conventional branch-and-bound solution is too slow for videos.

# 4.1.2 Spatio-temporal branch-and-bound search

We present a new method called spatio-temporal branchand-bound search (STBB) to search the video space. Instead of directly applying branch-and-bound in the 6D parameter space, our new method decomposes it into two subspaces: (1) 4D spatial parameter space and (2) 2D temporal parameter space. We denote by  $W \in \mathbb{R} \times \mathbb{R} \times \mathbb{R} \times \mathbb{R}$  a spatial window and  $T \in \mathbb{R} \times \mathbb{R}$  a temporal segment. A subvolume V is uniquely determined by W and T. The detection score of a subvolume  $f(V_{W \times T})$  is:

$$f(V_{W\times T}) = f(W,T) = \sum_{d \in W \times T} s(d).$$
(11)



Fig. 2. Illustration of the upper bound estimation.  $\mathbb{W}$  denotes a set of spatial windows.  $W_{max}$  and  $W_{min}$  are the maximum and minimum windows in  $\mathbb{W}$ , respectively. F and G are two matrices obtained through Eq. 16 and Eq 17, respectively. Empty cells in F and G matrices correspond to null entries. Left figure: the first upper bound in Lemma 1. The upper bound is  $\hat{F}_1(\mathbb{W}) = 19 + 9 + 7 = 35$ . Right figure: the second upper bound in Lemma 2. The upper bound is  $\hat{F}_2(\mathbb{W}) = 21 - (-3 - 9 - 1) = 34$ .

Let  $\mathbb{W} = [0, m] \times [0, m] \times [0, n] \times [0, n]$  be the parameter space of the spatial windows, and  $\mathbb{T} = [0, t] \times [0, t]$  be the parameter space of temporal segments. Our objective here is to find the spatio-temporal subvolume which has the maximum detection score:

$$W^*, T^*] = \arg \max_{W \in \mathbb{W}, T \in \mathbb{T}} f(W, T).$$
(12)

The optimal detection score is:

$$F(W^*) = \max_{W \in \mathbb{W}} F(W) = \max_{W \in \mathbb{W}} \max_{T \in \mathbb{T}} f(W, T).$$
(13)

We take different search strategies in the two subspaces  $\mathbb{W}$  and  $\mathbb{T}$  and search alternately between  $\mathbb{W}$  and  $\mathbb{T}$ . First, if the spatial window W is determined, we can easily search for the optimal temporal segment in space  $\mathbb{T}$ :

$$F(W) = \max_{T \in \mathbb{T}} f(W, \tilde{T}), \tag{14}$$

This relates to the max subvector problem, where given a real vector, the output is the contiguous subvector of the input that has the maximum sum (see Fig. 3). We will discuss its efficient solution later.

To search the spatial parameter space W, we employ a branch-and-bound strategy. Since the efficiency of a branch-and-bound based algorithm critically depends on the tightness of the upper bound, we first derive a tighter upper bound.

Given an arbitrary parameter space  $\mathbb{W} = [m_1, m_2] \times [m_1, m_2] \times [n_1, n_2] \times [n_1, n_2]$ , the optimal solution is:

$$W^* = \arg\max_{W \in \mathbb{W}} F(W). \tag{15}$$

We define  $F(\mathbb{W}) = F(W^*)$ . Assume there exist two subrectangles  $W_{min}$  and  $W_{max}$  such that  $W_{min} \subseteq W \subseteq$  $W_{max}$  for any  $W \in \mathbb{W}$ . For each pixel  $i \in W_{max}$ , we denote the maximum sum of the 1D subvector along the temporal direction at pixel *i*'s location by:

$$F(i) = \max_{T \subseteq \mathbb{T}} f(i, T).$$
(16)

Let  $F^+(i) = \max(F(i), 0)$ , we have the first upper bound for  $F(\mathbb{W})$ , as presented in Lemma 1.

Lemma 1: (upper bound  $\hat{F}_1(\mathbb{W})$ ) Given a spatial parameter space  $\mathbb{W} = \{W : W_{min} \subseteq W \subseteq W_{max}\}$ , we have

$$F(\mathbb{W}) \leq \hat{F}_1(\mathbb{W}) = F(W_{min}) + \sum_{i \in W_{max}, i \notin W_{min}} F^+(i).$$

When  $W_{max} = W_{min}$ , we have the tight bound  $\hat{F}_1(\mathbb{W}) = F(W_{min}) = F(W^*)$ .

Symmetrically, for each pixel  $i \in W_{max}$ , we denote the minimum sum of the 1D sub-vector at pixel *i*'s location by:

$$G(i) = \min_{T \subset \mathbb{T}} f(i, T).$$
(17)

Let  $G^{-}(i) = \min(G(i), 0)$ , and Lemma 2 presents the other upper bound of  $F(\mathbb{W})$ .

# *Lemma 2:* (upper bound $\hat{F}_2(\mathbb{W})$ )

Given a spatial parameter space  $\mathbb{W} = \{W : W_{min} \subseteq W \subseteq W_{max}\}$ , we have

$$F(\mathbb{W}) \le \hat{F}_2(\mathbb{W}) = F(W_{max}) - \sum_{i \in W_{max}, i \notin W_{min}} G^-(i).$$

When  $W_{max} = W_{min}$ , we have the tight bound  $\hat{F}_2(\mathbb{W}) = F(W_{max}) = F(W^*)$ .

The proofs of Lemma 1 and Lemma 2 are given in the Appendix. The two lemmas are illustrated in Fig. 2, where  $F(W_{min}) = 19$ , and  $F(W_{max}) = 21$ . The values of F(i) are shown in the F matrix where blank cells indicate zeros. The values of G(i) are shown in the G matrix. Lemma 1 gives the upper bound  $\hat{F}_1(\mathbb{W}) = 19 + 9 + 7 = 35$  and Lemma 2 gives the upper bound  $\hat{F}_2(\mathbb{W}) = 21 - (-3 - 9 - 1)$ 

1) = 34. Based on Lemma 1 and Lemma 2, we can obtain a final tighter upper bound, which is the minimum of the two available upper bounds.

Theorem 1: (Tighter upper bound  $\hat{F}(\mathbb{W})$ )

Given a spatial parameter space  $\mathbb{W} = \{W : W_{min} \subseteq W \subseteq W_{max}\}$ , the upper bound of the optimal solution  $F(\mathbb{W})$  is:

$$F(\mathbb{W}) \le \hat{F}(\mathbb{W}) = \min\{\hat{F}_1(\mathbb{W}), \hat{F}_2(\mathbb{W})\}.$$
 (18)

Based on the tighter upper bound derived in Theorem 1, we propose our new branch-and-bound solution in the spatial parameter space W in Alg. 2. Since the convergence speed of branch-and-bound method highly depends on the tightness of the bound, the new algorithm can converge much faster with a better upper bound estimation. Moreover, compared to the conventional branch-and-bound solution in Alg. 1, the new STBB algorithm keeps track of the current best solution which is denoted by  $W^*$ . Only when a parameter space W contains potentially better solution (*i.e.*  $\hat{F}(W) \geq F^*$ ), we push it into the queue. Otherwise, we discard the whole space W. It thus saves the memory in maintaining the priority queue of W.

Algorithm 2: Spatio-temporal branch-and-bound (STBB) search input : video  $\mathcal{V} \in \mathbb{R}^{m \times n \times t}$ quality bounding function  $\hat{F}$  (see text) **output** :  $V^* = \arg \max_{V \subseteq \mathcal{V}} f(V)$ 1 initialize P as an empty priority queue **2** set  $\mathbb{W} = [T, B, L, R] = [0, n] \times [0, n] \times [0, m] \times [0, m]$ 3 set  $\hat{F}(\mathbb{W}) = \min\{\hat{F}_1(\mathbb{W}), \hat{F}_2(\mathbb{W})\}$ 4 push ( $\mathbb{W}$ ,  $\hat{F}(\mathbb{W})$ ) into P 5 set current best solution  $\{W^*, F^*\} = \{W_{max}, F(W_{max})\};$ 6 repeat retrieve top state  $\mathbb{W}$  from P based on  $\hat{F}(\mathbb{W})$ 7 8 if  $(\hat{F}(\mathbb{W}) > F^*)$  then split  $\mathbb{W} \to \mathbb{W}^1 \cup \mathbb{W}^2$ 9 CheckToUpdate( $\mathbb{W}_1, W^*, F^*, P$ ); 10 CheckToUpdate( $\mathbb{W}_2, W^*, F^*, P$ ); 11 12 else  $T^* = \arg\max_{T \subseteq [0,t]} f(W^*, T);$ 13 return  $V^* = [W^{\overline{*}}, T^*]$ 14 15 until stop ; 16 Function CheckToUpdate( $\mathbb{W}, W^*, F^*, P$ ) 17 Get  $W_{min}$  and  $W_{max}$  of  $\mathbb{W}$ 

# 4.1.3 Efficient upper bound estimation for branchand-bound search

To estimate the upper bound in Theorem 1, as well as to search for the optimal temporal segment  $T^*$  given a



Fig. 3. Max subvector search: the highlighted element by red is the subvector of max sum, which is 8 - 1 + 5 = 12.

spatial window W, we design an efficient way to evaluate  $F(W_{max})$ ,  $F(W_{min})$ , and in general F(W).

According to Eq. 14, given a spatial window W of a fixed size, we need to search for a temporal segment with maximum summation. To present our efficient solution, we first review the classic max subvector problem in one-dimension pattern recognition. It is the degeneration of the maximum subvolume problem in spatio-temporal space. There exists an elegant solution called Kadane's algorithm which is of a linear complexity using dynamic programming [56]. We present Kadane's algorithm in Alg. 3. The max sum problem is illustrated in Fig. 3.

Kadane's algorithm can accelerate the temporal search and provide an efficient estimation of the upper bounds. Given any spatial window W, the summation within Wat each frame j is  $f(W, j) = \sum_{d \in W \times j} s(d)$ . By applying the trick of integral-image, f(W, j) can be obtained in a constant time. Let v(j) = f(W, j), the evaluation of F(W)in Eq. 14 is to find the max subvector in v. By using Kadane's algorithm, it can be done in a linear time. As a result, both upperbounds in Lemma 1 and Lemma 2 can be obtained in a linear time. Therefore the estimation of the upper bound  $\hat{F}(W)$  in Theorem 1 is of a linear complexity O(t).

The complexity comparison between our proposed method (Alg. 2) and the conventional branch-and-bound (Alg. 1) is presented in Table 1. As our branch-and-bound is only performed in the spatial space, the worst case complexity of our Alg. 2  $(O(m^2n^2t))$  is better than that of Alg. 1  $(O(m^2n^2t^2))$  which needs to perform branch-and-bound in the spatio-temporal space.

Algorithm 3: The linear algorithm of max subvec-									
tor [56]									
input : real vector v of length $t + 1$ output : $T^* = \arg \max_{T \subseteq [0,t]} \sum_{i \in T} v(i)$									
1 set $MaxSofar = MaxEndingHere = 0;$									
2 set $Start = End = 0;$									
<b>3 for</b> $i = 0: t$ <b>do</b>									
4 $MaxEndingHere =$									
$\max(0, MaxEndingHere + v(i));$									
5 if $MaxEndingHere = 0$ then									
$6 \qquad \qquad \mathbf{CurStart} = \min(i+1,t);$									
7 <b>if</b> $MaxSoFar \leq MaxEndingHere$ <b>then</b>									
8 $Start = CurStart;$									
9 $End = i;$									
10 $\[ MaxSofar = max(MaxSofar, MaxEndingHere); \]$									
11 return $T^* = [Start, End];$									

	<b>BB</b> (Alg. 1)	STBB (Alg. 2)
Dimensions for B&B	6 (spatio-temporal)	4 (spatial)
Upper bound est.	O(1)	O(t)
Worst case	$O(m^2n^2t^2)$	$O(m^2n^2t)$

Complexity comparison between spatio-temporal branch-and-bound and conventional branch-and-bound.

# 5 MULTI-CLASS MULTIPLE-INSTANCE AC-TION DETECTION

# 5.1 Multiple-instance detection algorithm

The branch-and-bound approaches in Alg. 1 and Alg. 2 are designed to search for a unique subvolume of maximum score. For multiple instance action detection, the same algorithm needs to be performed multiple rounds. At each round, the score of the detected subvolume  $V^*$  is compared against a predefined detection threshold  $D_t$  in order to determine whether it is a valid detection. If it is a valid detection, we clear it by setting the score of  $i \in V^*$  to  $s(d_{\emptyset})$  and continue to find the next subvolume of the maximum detection score. This process continues until the current best subvolume is not a valid detection.

Algorithm 4: Multiple-instance action detection	
input : video $\mathcal{V} \in \mathbb{R}^{m \times n \times t}$ ; detection threshold $D_t$ output : a collection of detections: $V^* \subseteq \mathcal{V}$ , s.t. $f(V^*) \ge D_t$	
1 repeat 2 $V^* = \text{STBBSearch}(\mathcal{V});$ 3 clear $V^*$ to zero values and update $\mathcal{V}.$ 4 until the current detection is invalid: $f(V^*) < D_t$ ;	

# 5.2 Accelerated STBB (A-STBB) for multipleinstance detection

To improve the efficiency of multiple-instance detection, we modify the original STBB search in Alg. 2 and propose an accelerated STBB (A-STBB). We briefly explain the main idea below. For multiple-instance detection, the detection threshold  $D_t > 0$  can be used to speed up the search process by terminating many unnecessary branches earlier during the branch-and-bound process. First of all, if there is no valid detection in a video sequence, then instead of finding the optimal subvolume  $V^*$  with the maximum detection score, we can safely terminate the search at an earlier stage. For example, if a parameter space  $\mathbb V$ satisfies  $\hat{f}(\mathbb{V}) \leq D_t$ , it indicates that  $\mathbb{V}$  is an invalid parameter space, because the score of the best candidate is still below the detection threshold. Therefore,  $\mathbb V$  does not require a further inspection. If none of the remaining candidates satisfies  $\hat{f}(\mathbb{V}) \geq D_t$ , then the search can be safely terminated because no valid detection will be found.

Furthermore, if a subvolume V with valid detection score  $f(V) \ge D_t$  is already found, we can quickly finalize the detection based on the current solution, instead of keeping looking for the maximum  $V^*$ . In such a case, although the final detection may not be the optimal subvolume  $V^*$ , it still provides a valid detection where  $f(V) \ge D_t$ . Therefore, it leads to a much faster search without significantly degrading the quality of the detection results.

Algorithm 5: Accelerated STBB (A-STBB) search
input : video $\mathcal{V} \in \mathbb{R}^{m \times n \times t}$ ; detection threshold $D_t$ output : a subvolume $\tilde{V} \subseteq \mathcal{V}$ , s.t. $f(\tilde{V}) \ge D_t$ (if no valid
detection, return $V = \emptyset$ )
1 set $\mathbb{W} = [T, B, L, R] = [0, n] \times [0, n] \times [0, m] \times [0, m]$ 2 get $\hat{F}(\mathbb{W}) = \min{\{\hat{F}_1(\mathbb{W}), \hat{F}_2(\mathbb{W})\}}$
3 push ( $\mathbb{W}$ , $\hat{F}(\mathbb{W})$ ) into empty priority queue P 4 set current best solution { $W^* F^*$ } = { $W_{max} F(W_{max})$ }:
5 repeat
6 retrieve ton state $\mathbb{W}$ from P based on $\hat{F}(\mathbb{W})$
7 if $\hat{F}(\mathbb{W}) < D$ , then
$\begin{array}{c c} & & I & I & (\forall V) \leq D_t \text{ then} \\ \hline & & & I \text{ return } \tilde{V} = \emptyset \end{array}$
9 <b>if</b> $(\hat{F}(\mathbb{W}) > F^*)$ then
10 split $\mathbb{W} \to \mathbb{W}^1 \cup \mathbb{W}^2$
11 CheckToUpdate( $\mathbb{W}_1, W^*, F^*, P$ );
12 Check To Update $(\mathbb{W}_2, W^*, F^*, P);$
13 else $T^* = \operatorname{ang} \max \left\{ \int W^* T \right\}$
$\begin{array}{c c} I & = \arg \max_{T \subset [0,t]} f(W, T), \\ return \tilde{V} & [W^*, T^*] \end{array}$
[ etuni v = [vv, I].
16 until stop ;
17 <b>FUNCTION</b> Check IOUpdate( $W, W^*, F^*, P$ ) 18 if $(P^* > D)$ then
18 II $(F \ge D_t)$ then to be a priority success $D$ such $(\mathbb{N}V \ \hat{E}(\mathbb{N}V))$ into constants
priority queue P push (w, F (w)) into empty
20 else

21 22	Get $W_{min}$ and $W_{max}$ of $\mathbb{W}$ if $(F(W_{min}) > F^*)$ then $\  \  \  \  \  \  \  \  \  \  \  \  \  $
23	if $(F(W_{max}) > F^*)$ then
24	update $\{W^*, F^*\} = \{W_{max}, F(W_{max})\};$
25	if $(W_{max} \neq W_{min})$ then
26	get $\hat{F}(\mathbb{W}) = \min\{\hat{F}_1(\mathbb{W}), \hat{F}_2(\mathbb{W})\}$
27	if $\hat{F}(\mathbb{W}) \geq F^*$ then
28	push ( $\mathbb{W}$ , $\hat{F}(\mathbb{W})$ ) into P

Incorporating the above two heuristics, we present the accelerated STBB (A-STBB) search in Alg. 5. Compared with the STBB search in Alg. 2, during each search iteration, we retrieve an upper bounded estimation  $\hat{F}(\mathbb{W})$  from the heap. If  $F(\mathbb{W}) < D_t$ , we directly reject the whole video sequence  $\mathcal{V}$ , since no  $V^*$  can achieve the detection threshold. This strategy largely speeds up the scanning of negative video sequences which do not contain the target action. Moreover, at each search iteration, we also keep track of the current best score  $F^*$ . When  $F^* \ge D_t$ , it indicates that there exists a valid detection in the corresponding parameter space  $\mathbb{W}$ . In such a case, we speed up the search by limiting the rest of the search space within  $\mathbb W$  only. In other words, instead of searching for the optimal  $f(V^*)$  globally, we are satisfied with the local optimal solution  $f(\tilde{V}) > D_t$ . Since only one subvolume with qualified score will be selected while other subvolumes are discarded, our A-STBB performs the nonmaxima suppression implicitly during the search process.

# 6 EXPERIMENTS

#### 6.1 Action Categorization

We use the KTH dataset to evaluate the proposed NBMIM classifier on action categorization. The KTH dataset contains six types of human actions: walking, jogging, running, boxing, hand waving and hand clapping, each of which is performed several times by 25 subjects. There are 4 different environments where the video sequences are captured: outdoors, outdoors with scale variation, outdoors with different clothes and indoors. The video is captured at 25 frames per second and at a low image resolution of  $160 \times 120$ .

We follow the standard experiment setting of KTH dataset as in [36] [19]. The whole dataset contains 598 video sequences, taken over homogeneous backgrounds with a static camera. Each sequence is further segmented into 4 subsequences according to [36], thus it gives in total 2391 action videos. Each action video has an average length of four seconds. Among the 25 persons, 16 of them are used for training and the rest 9 are used for testing. The training dataset contains 1528 individual actions and the testing dataset contains 863 individual actions. We apply both motion (histogram of motion) and appearance (histogram of gradient) features as in [19]. By concatenating the HOG and HOF features, a 162-dimensional feature vector is used to characterize each STIP. The average Euclidean length of the STIP descriptor is 4.46. The training dataset generates a pool of 308, 110 STIPs. Given a query STIP, we search its  $\epsilon$ -nearest neighbors using locality sensitive hashing. The E2LSH package [55] is employed and the probability for correct retrieval is set to p = 0.9.

The threshold of the nearest neighbor search,  $\epsilon$ , is the only parameter of the proposed NBMIM classifier. Its influence is two-fold. First of all, it affects the search speed and the quantity of the nearest neighbors. The larger the  $\epsilon$ , the slower the approximate  $\epsilon$ -NN search using LSH, but the more nearest neighbors it will find. Secondly,  $\epsilon$  also controls the bandwidth  $\sigma$  in the kernel density estimation according to Eq. 6. To evaluate the influence of  $\epsilon$ , we test different choices of  $\epsilon$  and compare three different classification models: NBMIM (adaptive kernel bandwidth), NBMIM (fixed kernel bandwidth), and NBNN in [52]. To make a fair comparison to NBNN, we use the same parameter for the approximate nearest neighbor search as described in Section 3.3.1. All of the three classifiers share the same  $d_{NN}^{c+}$  and  $d_{NN}^{c-}$ . The only difference is the voting score  $s_c(d)$ . In this experiment, since each action class has approximately the same number of video sequences, we assume the prior probabilities are equal and apply Eq. 3 to calculate  $s_c(d)$ . The result in Table 2 shows that the classification performance is not very sensitive to the selection of  $\epsilon$ . Our proposed NBMIM with the adaptive kernel bandwidth performs slightly better than NBMIM with a fixed bandwidth, as well as NBNN. It is worth noting

that the NBNN classifier cannot be directly applied to the detection formulation of Eq. 7, because its voting score is always positive.

	$\epsilon = 1.8$	$\epsilon = 2.0$	$\epsilon = 2.2$	$\epsilon = 2.4$	$\epsilon = 2.6$					
adaptive	91.8%	93.0%	93.7 %	93.4 %	93.3 %					
$\frac{1}{2\sigma^2} = 1$	91.9%	92.2%	92.7 %	92.7 %	92.6 %					
NBNN	91.7%	91.8%	92.5%	92.6%	92.7 %					

TABLE 2

The comparison between NBMIM (adaptive and fixed kernel bandwidth) and NBNN, with different selections of  $\epsilon$ .

The best action categorization results are presented in Table 3, with  $\epsilon = 2.2$  and using the adaptive kernel bandwidth. Among the 863 testing actions, we obtained 54 errors, and the total accuracy is 93.7%. Among the six types of actions, hand clapping, walking and boxing receive 100% accuracy. Most of the errors are due to the mis-classification of running to jogging.

	clap	wave	walk	box	run	jog						
clapping	144	0	0	0	0	0						
waving	5	139	0	0	0	0						
walking	0 0		144	0	0	0						
boxing	0	0	0	143	0	0						
running	1	0	0	0	105	38						
jogging	2	0	4	0	4	134						

TABLE 3 Confusion matrix for the KTH action dataset. The total accuracy is 93.7%.

In Table 4, we further compare our results with that of [19], by applying exactly the same training <sup>1</sup> and testing dataset, as well as the same STIP features. However, we do not quantize STIPs into "words". Instead of using the SVM, we match the raw STIPs in the original highdimensional feature space and apply the NBMIM classifier. Our results show that, without quantizing primitive features into "words", the classification performance can be further improved. This is consistent with the discussion in [52] which pointed out that the nearest neighbor approach has the potential to provide better classification performance than the SVM based on the "bag-of-words" representation, where the quantization step can introduce a loss of discriminative information.

#### 6.2 Detecting two-hand waving action

We select the two-hand waving action as a concrete example for action detection. To validate the generalization ability of our method, we apply completely different datasets for training (KTH dataset) and testing (CMU action dataset [4]). As summarized in Table 5, for the positive training data, we apply the KTH hand waving dataset that contains 16 persons. The negative training data is constituted by two parts (1) the KTH walking dataset which contains 16

<sup>1.</sup> In [19], 8 persons are used for training and another 8 persons are used as cross-validation for parameter tuning of the SVM. We use the whole 16 persons as the training data.

	training	training testing features				
[19]	8 persons + 8 person CV	9 persons	STIP + "bag of words"	non-linear SVM	91.8 %	
ours	16 persons	9 persons	STIP	NBMIM	93.7 %	

TABLE 4 Comparison between NBMIM and SVM.



Fig. 4. Performance of two-hand wave detection in CMU dataset, with different choices of  $s(d_{\emptyset})$ . See text for the measurement of precision and recall scores.

persons and (2) one office indoor sequence which contains actions of sitting down and standing up.

The testing dataset has 48 sequences, which include two types of actions in the CMU dataset: (1) two-hand waving and (2) jumping jacks. Both of them contain the two-hand waving actions. The duration of each test sequence ranges from 10 to 40 seconds, and the video resolution is  $160 \times 120$ . Among the 48 sequences, 19 of them contain a total number of 52 positive instances. The other 29 sequences do not contain positive examples. For the density estimation, we set  $\epsilon = 2.6$  for the nearest neighbor search. To avoid noisy detection, we regard a detection as invalid if its temporal duration is shorter than 0.4 seconds.

To evaluate the results, we apply a similar measurement proposed in [4], but with a relatively loose criterion. For the precision score, a detection is regarded as correct if at least 1/8 of the volume size overlaps with the ground truth. For the recall score, the ground truth is regarded as retrieved if at least 1/8 of its volume size is covered by at least one detection. We use such a measurement because the ground truth labels cover the entire human body while the region that exhibits hand waving action is much smaller. Even though the localized bounding box of the proposed method is not very accurate, one can always refine the bounding box by using other cues after roughly determining the location. We apply the precision and recall scores to evaluate the detection performance, where precision = # correct detect / # total detect, and recall = # correct detect / # total action.

positive training	hand-waving 16 persons (KTH)
negative training	walking 16 persons (KTH) + 1 indoor seq.
testing	two-hand waving + jumping jacks (CMU)

TABLE 5 Cross-dataset training and testing of two-hand waving detection.

We apply the efficient A-STBB search to detect two hand waving actions. To evaluate the influence of the parameter  $s(d_{\emptyset})$ , we test a number of different values of  $s(d_{\emptyset})$ , including  $s(d_{\emptyset}) = -10 \times 10^{-5}$ ,  $-7.5 \times 10^{-5}$ ,  $-6 \times 10^{-5}, -5 \times 10^{-5}, -4 \times 10^{-5}, -2.5 \times 10^{-5}, -1 \times 10^{-5}.$ Fig. 4 presents the precision-recall curves by increasing the detection threshold  $D_t$  from 5 to 40. It shows that  $s(d_{\emptyset})$ is an important parameter that can influence the detection results significantly. When a small  $s(d_{\emptyset})$  is selected, the detected maximum subvolume is of a large size, thus having a sufficient overlap with the ground truth. Therefore, we obtain a higher recall score while the precision score gets worse. On the other hand, when a large  $s(d_{\emptyset})$  is selected, the detected subvolume is of a small size thus the overlap with the ground truth volume becomes smaller. This results in a worse recall score but a better precision score. When selecting  $s(d_{\emptyset}) = -4 \times 10^{-5}$ , both precision and recall scores achieve above 70% at a specific detection threshold.

Some detection examples are presented from Fig. 10 to Fig. 14. The yellow bounding box is the ground truth label of the whole human body action and the red bounding box is our detection of the two-hand waving action. Since both motion and appearance features are used, our method can tolerate action pattern variations caused by the change of subjects. Our detection method can also handle scale changes of the actions, performing speed variations, background clutter, and even partial occlusion. Fig. 10 shows the same person performing two-hand waving with two different styles and different speeds. In Fig. 11, two actions with large scale variations are detected successfully. Fig. 12 shows action detection results on cluttered backgrounds and with severe partial occlusions, where target tracking is very difficult.

Most of the missing and false detections are caused by the bad lighting conditions, crowded scenes, large view point changes, or moving cameras. Fig. 13 presents a false detection example where two single-hand waving actions occur together and brings a false detection of the twohand waving action. As the naive-Bayes assumption does not consider the geometric relations among STIPs, our approach cannot distinguish whether the two waving hands are from the same person or not. To better handle this problem, a geometric model would be required for a further verification. Finally, Fig. 14 shows an example of the missed detection. Although it is indeed a partial detection, the overlap region with the ground truth is less than 1/8, thus it is treated as a missed detection.

#### 6.3 Multi-class multiple-instance action detection

Based on the multi-class action recognition model, we can perform multi-class multiple-instance action detection.



Fig. 5. Examples of multi-class multiple-instance action detection. Each row shows an action class: hand waving (first row), hand clapping (second row) and boxing (third row). Each image is a sample frame from the action video. The first column shows the training videos from the KTH dataset. The second to the fifth columns show some detection results. The highlighted bounding boxes correspond to the detected spatial window  $W^*$  and we apply different colors to distinguish different action classes: clapping (turquoise), waving (magenta), and boxing (yellow). The final column shows miss detection examples, where the bounding boxes are ground truth labeling: clapping (red), waving (green), and boxing (blue). The whole dataset is accessible at http://research.microsoft.com/en-us/downloads/fbf24c35-a93e-4d22-a5fe-bc08f1c3315e/.

To validate the generalization ability of our method, we still apply a cross-dataset training and testing. We select three classes of actions for positive training from the KTH dataset: boxing, hand waving and hand clapping, including 16 subjects for each class of action. Because the action instances are captured in different environments and view points, and exhibit spatial scale and style variations, the intra-class variations of actions are well captured in the training data. To better distinguish the three types of target actions from other types of movements, we also use the walking class in the KTH dataset as the common negative class. As a result, for each of the three action classes, the negative training dataset includes the STIPs from the walking class, as well as the STIPs from other two action classes.

The testing videos are captured by ourselves. Each testing sequence is of a higher resolution  $320 \times 240$ , compared with that of  $160 \times 120$  in the training videos in the KTH dataset. The frame rate is 15 frames per second. The testing dataset contains 16 video sequences. Each video sequence is between 32 and 76 seconds. It has in total 63 action instances: 14 hand clapping, 24 hand waving, and 25 boxing, performed by 10 different subjects who do not appear in the training data. Each sequence contains multiple types of actions, performed by one or multiple subjects. As a challenging dataset, all of the video sequences are captured in cluttered and moving backgrounds, including both indoor and outdoor scenes. The style and scale of actions can vary significantly depending on the subject. To evaluate the performance, we manually label a spatiotemporal bounding box for each action instance. A detected action is regarded as correct if at least 1/8 of the volume size overlaps with a ground truth label. On the other hand, an action is regarded as retrieved if at least 1/8 of its volume size overlaps with that of a valid detection. To filter out noisy detections, we require a valid detection lasts between 20 and 200 frames. For the kernel density estimation, we set the nearest neighbor search parameter to be  $\epsilon = 2.6$ .

We apply the A-STBB search for multi-class multipleinstance action detection. In Fig. 6, we show the precision and recall curves for three action classes, by increasing the detection threshold  $D_t$  from 3 to 30. We also compare a few different values of  $s(d_{\emptyset})$ , including  $s(d_{\emptyset}) = -1 \times 10^{-5}, -2 \times 10^{-5} - 3 \times 10^{-5}, -4 \times 10^{-5},$ and  $-6 \times 10^{-5}$ . For different action classes, the optimal parameters of  $D_t$  and  $s(d_{\emptyset})$  may be different. Among the three classes of actions, hand waving and boxing provide better performance, where both precision and recall rates are higher than or close to 65%. However, hand clapping is more challenging, especially if the clapping movement is subtle. Hand clapping is also easily confused with the hand waving action. For all of the three classes, most missing detections are due to the small spatial scales, bad lighting conditions, or crowded scences. In Fig. 7, we show the detection results of multiple actions in the same scene.

The computational cost of multi-class multiple-instance action detection contains three parts: (1) extraction of STIPs; (2) kernel density estimation and calculation of voting scores for each class; and (3) search for qualified subvolumes for each class. First, for videos at resolution  $320 \times 240$ , the speed of STIP detection is 2-4 frames per second using the binary code provided by [19]. Second, the major cost of obtaining the voting score  $s_c(d)$  comes from the  $\epsilon$ -NN search in density estimation. By using the E2LSH code for efficient NN search, the query time of each STIP is 40-50 milliseconds with  $\epsilon = 2.6$  and retrieve probability p = 0.9. However, if performing exhaustive search of  $\epsilon$ -NN, the query time of each STIP increases to 130 milliseconds. If parallel search can be performed using a four-core CPU, the estimated query time can achieve around 12 milliseconds per STIP. As each frame



Fig. 6. Performance of three-class action detection. Left: hand clapping; middle: hand waving; right: boxing.



Fig. 7. Detection of multiple actions in the same scene. Top: detection of two boxings (yellow). Bottom: detection of one boxing (yellow) and one hand waving (purple).

contains 20-40 STIPs on average, the processing time can achieve 2-4 frames per second. Finally, to evaluate the CPU cost of subvolume search through A-STBB, we record the computational cost for each of the 16 testing sequences in Table 6. The test is performed on a four-core CPU desktop.

In Table 6, we notice that the computational cost of A-STBB depends on the video sequence, including the number of STIPs and the number of action instances. On average, the A-STBB search can achieve 4-5 frames per second using a four-core CPU. The search tends to be slower for video sequences with a larger number of moving objects in the background since a lot of STIPs will be extracted. On the other hand, if a video sequence does not contain any target actions, the search will finish quickly thanks to the early termination strategy.

# 6.4 Search complexity and efficiency comparison

# 6.4.1 Comparison between STBB search and conventional BB search

To validate the efficiency gain of our STBB search, we compare our STBB (Alg. 2) to the conventional branchand-bound (Alg. 1) by searching the MVI-142a sequence in the CMU action dataset [4]. The max subvolume is of size  $43 \times 32 \times 112$ . The input video  $\mathcal{V}$  is of size  $120 \times 160 \times 141$ , a temporal segment from MVI-142a. We intentionally choose such a target video of a short length, such that the sizes of its 3 dimensions are balanced. This gives a fair comparison to the conventional branch-and-bound, because the longer the video length *t*, the less efficient the conventional branch-and-bound will be.

The left figure in Fig. 8 shows that our proposed method converges much faster than the conventional branch-andbound. In terms of the number of branches, our method converges after 10,302 branches, an order of magnitude faster than the conventional branch-and-bound which needs 103,202 branches before convergence. This validates that the upper bound proposed in Theorem 1 is tighter than that of the conventional method. In Alg. 1, the upper bounded estimation  $\hat{f}(\mathbb{V})$  decreases slowly when the current state converges to the optimal solution. In comparison, the convergence of the upper bound in our proposed method (Alg. 2) is much faster. For example, after 2000 branches, our method reaches a very good solution f(V) = 15.78, which is close to the optimal one  $f(V^*) = 16.21$ . On the other hand, after 2000 branches, the largest upper bound given by the conventional branch-and-bound is still as large as  $\hat{f}(\mathbb{V}) = 24.06$ .

As mentioned earlier, another advantage of our method is that it keeps track of the current best solution. A new subvolume is pushed into the queue only when its upper bound is better than the current best solution. In comparison, the method proposed in [1] needs to push every middle state into the priority queue, as there is no record of the current best solution. In Fig. 8, we also compare the required size of the priority queue between our method and the conventional branch-and-bound. The size of the priority queue in our method is well controlled and is much smaller. In our method, during the branch-and-bound process, the size of the priority queue decreases after a peak value. However, for the conventional branch-and-bound, the size of priority queue always increases, almost linearly to the number of branches. Since each insertion or extraction operation of priority queue is O(logn) for a queue of size n, the size of the priority queue affects both the computational and memory costs. It is especially important to limit a queue to a moderate size for the video space search because it can generate a much larger number of candidates than the spatial image case.

# 6.4.2 Evaluation of the accelerated STBB search

To evaluate the efficiency of the proposed accelerated STBB (A-STBB) algorithm in Alg. 5 for branch-andbound, we select the first five video sequences of the two-hand waving action in the CMU action dataset [4]. Each sequence contains one two-hand waving action. The

video #	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16
length #min	0.75	0.73	0.95	0.57	1	0.82	0.65	0.85	0.85	1.05	0.90	0.63	0.53	1.27	0.72	0.88
cost #min	0.25	2	3	1.5	1.5	6	0.75	0.25	2.75	0.5	1.25	0.25	0.25	10	3.75	10.25

#### TABLE 6

The CPU cost of the A-STBB search on the MSR action dataset for multiple-instance detection of three classes of actions. The algorithm is implemented in C++ and runs on a four-core CPU desktop. Only the A-STBB search cost is listed, while the STIP extractions and score calculation are not included. The measurement of CPU cost is rounded into  $\frac{1}{4}$  minutes.



Fig. 8. Comparison between our STBB search (Alg. 2) and conventional BB search in 3-dimensional spatio-temporal space (Alg. 1). Left: Comparison of the convergence speed. Both methods can find the optimal subvolume with detection score  $f(V^*) = 16.21$ . The red-circle curve and the blue-square curve show the convergence of the upper bound estimation of STBB method and conventional BB. Right: comparison between the length of the priority queue in branch-and-bound search.

algorithm searches for the subvolume of high detection score, such that it covers the action. Compared with original STBB which targets the optimal subvolume with maximum score, A-STBB finds approximate optimal subvolume but at a faster search speed. For the original STBB search, we do not need to specify the detection threshold, as it returns the subvolume with maximum detection score. For the accelerated STBB (A-STBB) search, the detection threshold is selected as  $D_t = 10$ . Under this detection score, the first subvolume returned by A-STBB is compared with the optimal subvolume returned from the original STBB algorithm. As the detection score of all of the five target subvolumes is higher than  $D_t = 10$ , such a detection threshold will not affect the efficiency comparison between the original STBB and A-STBB search algorithms.

The comparison between the A-STBB in Alg. 5 with the original STBB in Alg. 2 is presented in Table 7.  $W^*$  is the spatial window containing left, right, top, and bottom parameters.  $T^*$  includes the start and end frames. Table 7 shows that detection results of A-STBB in Alg. 5 are close to those of STBB in Alg. 2. Both algorithms provide similar detections results, in terms of detection scores, locations and sizes of the subvolumes. However, the number of branches in STBB can be up to 20 times more than that of A-STBB. It validates the efficiency of the proposed A-STBB. Moreover, if a video sequence does not contain any target actions, A-STBB can be even more efficient by terminating the search process at a very early stage, and returns with non-valid detection found.

To show the performance of our A-STBB search in real video surveillance scenario, we also test on two sequences from the TRECVID 2008 event detection dataset, which is a very challenging one for video surveillance [57]. The videos are captured by the real surveillance cameras in an airport. Although there are a lot of actions defined in

TRECVID 2008, we only use the running action since it is similar to those in the KTH dataset. We use 16 running persons from the KTH dataset for positive training and 16 walking persons for negative training. The two selected testing sequences are taken by the 2nd camera (five cameras in total). The video resolution is  $180 \times 144$  with 25 frames per second. Fig. 9 shows the detection results where each row corresponds to a video sequence.

# 7 CONCLUSION

Similar to the sliding window based search for object detection, detection of actions is to search for qualified subvolumes in the volumetric video space. To address the search complexity of this new formulation of action detection, a novel spatio-temporal branch-and-bound (STBB) search solution is proposed. We extend the previous branch-and-bound solution from searching spatial image patterns to searching spatio-temporal video patterns. By tightening the upper bound and reducing the parameter space from 6 dimensions to 4 dimensions, the STBB search is significantly more efficient in searching video patterns. For multi-class multiple-instance action detection, the accelerated STBB (A-STBB) search validates its efficiency and effectiveness on the CMU and MSR datasets.

In order to tolerate the intra-class action variations, we propose a discriminative pattern matching method, called naive-Bayes mutual information maximization (NBMIM), for action classification. Compared with conventional template-based pattern matching, instead of using a single template for pattern matching, we apply both positive and negative templates for discriminative matching. Despite its simplicity, the proposed NBMIM approach can well distinguish one action class from other classes, as well as the background class. Although such a naive-Bayes

video	$W^*$	$T^*$	score $f(V^*)$	# of branches
V1: STBB (Alg. 2)	55 97 23 54	442 553	16.21	206933
V1: A-STBB (Alg. 5)	55 97 23 52	442 546	15.97	4549
V2: STBB (Alg. 2)	61 122 20 38	673 858	37.39	67281
V2: A-STBB (Alg. 5)	60 122 20 39	673 858	37.36	4668
V3: STBB (Alg. 2)	72 118 22 71	11 700	89.01	71594
<b>V3</b> : A-STBB (Alg. 5)	82 114 23 73	10 705	85.21	2275
V4: STBB (Alg. 2)	73 112 23 77	420 1083	73.42	63076
V4: A-STBB (Alg. 5)	77 108 23 78	420 1083	70.93	2363
V5: STBB (Alg. 2)	18 144 7 114	418 451	46.50	315770
<b>V5</b> : A-STBB (Alg. 5)	41 151 7 114	419 451	45.36	133027

## TABLE 7

Comparison between STBB in Alg. 2 and accelerated STBB (A-STBB) in Alg. 5.  $V^* = [W^*, T^*]$  is the detected subvolume through branch-and-bound search. The four parameters of  $W^*$  determine the spatial location and the two parameters of  $T^*$ determine the start and end frames.



Fig. 9. Detection of person running in surveillance video using A-STBB search. Each row shows a detection instance. Top: a man runs in the corridor. Bottom: a child runs in the corridor.

assumption ignores the spatio-temporal dependency among interest point features, it leads to a better tolerance of intra-pattern variations. Our action detection method does not rely on the detection and tracking of a person. It can well handle scale changes, performing speed and style variations of actions, cluttered and dynamic backgrounds, even partial occlusions. The future work includes extending the STBB search to find subvolumes of more flexible shapes, *i.e.*, non-rectangle shapes, and relaxing the naive-Bayes assumption in discriminative matching to consider the spatio-temporal dependency among the interest points.

#### ACKNOWLEDGMENT

This work was supported in part by the Nanyang Assistant Professorship to Dr. Junsong Yuan, the National Science Foundation grant IIS-0347877, IIS-0916607, and US Army Research Laboratory and the US Army Research Office under grant ARO W911NF-08-1-0504. We thank Dr. Yan Ke, Dr. Cha Zhang and Dr. Zhengyou Zhang for helpful discussions, and Liangliang Cao for the help on the experiments of the TRECVID dataset.

# APPENDIX A

We prove the upper bound in Lemma 1 here.

$$f(W^*, T^*) = F(W^*) = \sum_{i \in W^* \times T^*} s(i)$$
  
$$= \sum_{i \in W_{min} \times T^*} s(i) + \sum_{i \in (W^* \setminus W_{min}) \times T^*} s(i)$$
  
$$\leq F(W_{min}) + \sum_{i \in (W^* \setminus W_{min})} F(i)$$
  
$$\leq F(W_{min}) + \sum_{i \in (W^* \setminus W_{min})} F^+(i)$$
  
$$\leq F(W_{min}) + \sum_{i \in (W_{max} \setminus W_{min})} F^+(i),$$

where  $i \in (W_1 \setminus W_2)$  denotes  $i \in W_1, i \notin W_2$ . When  $W_{max} =$  $W_{min}$ , we have  $\sum_{i \in (W_{max} \setminus W_{min})} F^+(i) = 0$ , which gives the tight bound  $F(W^*) = F(W_{min})$ .

# APPENDIX B

We prove the upper bound in Lemma 2 here.

T (TT 7\*)

$$f(W^*, T^*) = F(W^*) = \sum_{i \in W^* \times T^*} s(i)$$

$$= \sum_{i \in W_{max} \times T^*} s(i) - \sum_{i \in (W_{max} \setminus W^*) \times T^*} s(i)$$

$$\leq F(W_{max}) - \sum_{i \in (W_{max} \setminus W^*)} G(i)$$

$$\leq F(W_{max}) - \sum_{i \in (W_{max} \setminus W^*)} G^-(i)$$

$$\leq F(W_{max}) - \sum_{i \in (W_{max} \setminus W_{min})} G^-(i).$$

When  $W_{max} = W_{min}$ , we have  $\sum_{i \in (W_{max} \setminus W_{min})} G^{-}(i) = 0$ , which gives the tight bound  $F(W^*) = F(W_{max})$ .

# REFERENCES

- [1] C. H. Lampert, M. B. Blaschko, and T. Hofmann, "Beyond sliding windows: Object localization by efficient subwindow search," in Proc. IEEE Conf. on Computer Vision and Pattern Recognition, 2008
- [2] P. Viola and M. J. Jones, "Robust real-time face detection," Intl. Journal of Computer Vision, vol. 57, no. 2, pp. 137-154, 2004.
- [3] M. B. Blaschko and C. H. Lampert, "Learning to localize objects with structured output regression," in Proc. European Conf. on Computer Vision, 2008.
- [4] Y. Ke, R. Sukthankar, and M. Hebert, "Event detection in crowded videos," in Proc. IEEE International Conf. on Computer Vision, 2007.
- A. F. Bobick and J. W. Davis, "The recognition of human movement [5] using temporal templates," IEEE Trans. on Pattern Analysis and Machine Intelligence (PAMI), vol. 23, no. 3, pp. 257-267, 2001.

- [6] I. Laptev, "On space-time interest points," *Intl. Journal of Computer Vision*, vol. 64, no. 2-3, pp. 107–123, 2005.
- [7] C. Rao, A. Yilmaz, and M. Shah, "View-invariant representation and recognition of actions," *International Journal of Computer Vision* (*IJCV*), vol. 50, no. 2, pp. 203–226, 2002.
- [8] N. Nguyen, D. Phung, S. Venkatesh, and H. Bui, "Learning and detecting activities from movement trajectories using the hierarchical hidden markov models," in *Proc. IEEE Conf. on Computer Vision* and Pattern Recognition, 2005.
- [9] S. Ali, A. Basharat, and M. Shah, "Chaotic invariants for human action recognition," in *Proc. IEEE International Conf. on Computer Vision*, 2007.
- [10] V. Parameswaran and R. Chellappa, "View invariance for human action recognition," *International Journal of Computer Vision (IJCV)*, vol. 66, no. 1, pp. 83–101, 2006.
- [11] J. Sun, X. Wu, S. Yan, L. Cheong, T. Chua, and J. Li, "Hierarchical spatio-temporal context modeling for action recognition," in *Proc. IEEE Conf. on Computer Vision and Pattern Recognition*, 2009.
- [12] F. Lv and R. Nevatia, "Single view human action recognition using key pose matching and viterbi path searching," in *Proc. IEEE Conf.* on Computer Vision and Pattern Recognition, 2007.
- [13] D. Weinland and E. Boyer, "Action recognition using exemplar-based embedding," in *Proc. IEEE Conf. on Computer Vision and Pattern Recognition*, 2008.
- [14] A. Yilmaz and M. Shah, "Actions as objects: a novel action representation," in *Proc. IEEE Conf. on Computer Vision and Pattern Recognition*, 2005.
- [15] M. Blank, L. Gorelick, E. Shechtman, M. Irani, and R. Basri, "Actions as space-time shapes," in *Proc. IEEE International Conf.* on Computer Vision, 2005.
- [16] P. Dollar, V. Rabaud, G. Cottrell, and S. Belongie, "Behavior recognition via sparse spatio-temporal features," in *Proc. IEEE Intl. Workshop Visual Surveillance Performance Evaluation Tracking Surveillance*, pp. 65 – 72, 2005.
- [17] J. Liu, J. Luo, and M. Shah, "Recognizing realistic actions from videos "in the wild"," in *Proc. IEEE Conf. on Computer Vision and Pattern Recognition*, 2009.
- [18] J. Liu and M. Shah, "Learning human actions via information maximization," in *Proc. IEEE Conf. on Computer Vision and Pattern Recognition*, 2008.
- [19] I. Laptev, M. Marszalek, C. Schmid, and B. Rozenfeld, "Learning realistic human actions from movies," in *Proc. IEEE Conf. on Computer Vision and Pattern Recognition*, 2008.
- [20] K. Jia and D.-Y. Yeung, "Human action recognition using local spatio-temporal discriminant embedding," in *Proc. IEEE Conf. on Computer Vision and Pattern Recognition*, 2008.
- [21] J. Niebles, H. Wang, and L. Fei-Fei, "Unsupervised learning of human action categories using spatial-temporal words," *International Journal of Computer Vision*, vol. 79, no. 3, pp. 299–318, 2008.
- [22] I. Laptev, B. Caputo, C. Schu, and T. Lindeberg, "Local velocityadapted motion events for spatio-temporal recognition," *Computer Vision and Image Understanding*, vol. 109, no. 1, pp. 207–229, 2007.
- [23] P. S. Dhillon, S. Nowozin, and C. H. Lampert, *Combining Appearance and Motion for Human Action Classification in Videos*. Technical Report, Max-Planck-Institute for Biological Cybernetics, 2008.
- [24] P. Scovanner, S. Ali, and M. Shah, "A 3-dimensional sift descriptor and its application to action recognition," in *Proc. ACM Multimedia*, 2007.
- [25] Y. Wang and G. Mori, "Human action recognition by semi-latent topic models," *IEEE Trans. on Pattern Analysis and Machine Intelligence*, vol. 31, no. 10, pp. 1762–1774, 2009.
- [26] S. Ali and M. Shah, "Human action recognition in videos using kinematic features and multiple instance learning," *IEEE Trans. on Pattern Analysis and Machine Intelligence*, vol. 32, no. 2, pp. 288– 303, 2010.
- [27] A. Fathi and G. Mori, "Action recognition by learning mid-level motion features," in *Proc. IEEE Conf. on Computer Vision and Pattern Recognition*, 2008.
- [28] A. A. Efros, A. C.Berg, G. Mori, and J. Malik, "Recognizing action at a distance," in *Proc. IEEE International Conf. on Computer Vision*, 2003.
- [29] Z. Zhang, Y. Hu, S. Chan, and L.-T. Chia, "Motion context: A new representation for human action recognition," in *Proc. European Conf. on Computer Vision*, 2008.

- [30] P. Natarajan and R. Nevatia, "View and scale invariant action recognition using multiview shape-flow models," in *Proc. IEEE Conf. on Computer Vision and Pattern Recognition*, 2008.
- [31] S. N. Vitaladevuni, V. Kellokumpu, and L. S. Davis, "Action recognition using ballistic dynamics," in *Proc. IEEE Conf. on Computer Vision and Pattern Recognition*, 2008.
- [32] Z. Lin, Z. Jiang, and L. S. Davis, "Recognizing actions by shapemotion prototype trees," in *Proc. IEEE Intl. Conf. on Computer Vision*, 2009.
- [33] J. Liu, S. Ali, and M. Shah, "Recognizing human actions using multiple features," in *Proc. IEEE Conf. on Computer Vision and Pattern Recognition*, 2008.
- [34] D. Han, L. Bo, and C. Sminchisescu, "Selection and context for action recognition," in *Proc. IEEE Intl. Conf. on Computer Vision*, 2009.
- [35] S.-F. Wong, T.-K. Kim, and R. Cipolla, "Learning motion categories using both semantic and structural information," in *Proc. IEEE Conf.* on Computer Vision and Pattern Recognition, 2007.
- [36] C. Schüldt, I. Laptev, and B. Caputo, "Recognizing human actions: A local svm approach," in *ICPR'04*, vol. 3, (Cambridge, UK), pp. 32– 36, Aug. 23-26, 2004.
- [37] K. K. Reddy, J. Liu, and M. Shah, "Incremental action recognition using feature-tree," in *Proc. IEEE Intl. Conf. on Computer Vision*, 2009.
- [38] Y. Ke, R. Sukthankar, and M. Hebert, "Efficient visual event detection using volumetric features," in *Proc. IEEE International Conf.* on Computer Vision, 2005.
- [39] E. Shechtman and M. Irani, "Space-time behavior based correlation," in Proc. IEEE Conf. on Computer Vision and Pattern Recognition, 2005.
- [40] J. Yuan, Z. Liu, and Y. Wu, "Discriminative subvolume search for efficient action detection," in *Proc. IEEE Conf. on Computer Vision* and Pattern Recognition, 2009.
- [41] Y. Hu, L. Cao, F. Lv, S. Yan, Y. Gong, and T. S. Huang, "Action detection in complex scenes with spatial and temporal ambiguities," in *Proc. IEEE Intl. Conf. on Computer Vision*, 2009.
- [42] J. Yuan and Z. Liu, "Video-based human action detection sources," *IEEE Signal Processing Magazine*, vol. 27, no. 5, pp. 136 – 139, 2010.
- [43] C. Yeo, P. Ahammad, K. Ramchandran, and S. S. Sastry, "High-speed action recognition and localization in compressed domain videos," *IEEE Trans. on Circuits and Systems for Video Technology*, vol. 18, no. 8, pp. 1006–1015, 2008.
- [44] J. Yuan, Z. Liu, Y. Wu, and Z. Zhang, "Speeding up spatiotemporal sliding-window search for efficient event detection in crowded videos," in *Proc. ACM Multimedia Workshop on Events* in *Multimedia*, 2009.
- [45] M. D. Rodriguez, J. Ahmed, and M. Shah, "Action MACH a spatio-temporal maximum average correlation height filter for action recognition," in *Proc. IEEE Conf. on Computer Vision and Pattern Recognition*, 2008.
- [46] D. Weinland, R. Ronfard, and E. Boyer, "Free viewpoint action recognition using motion history volumes," *Computer Vision and Image Understanding*, vol. 104, no. 2-3, pp. 207–229, 2006.
- [47] H. Jiang, M. S. Drew, and Z.-N. Li, "Successive convex matching for action detection," in *Proc. IEEE Conf. on Computer Vision and Pattern Recognition*, 2006.
- [48] W. Li, Z. Zhang, and Z. Liu, "Expandable data-driven graphical modeling of human actions based on salient postures," *IEEE Trans.* on Circuits and Systems for Video Technology, vol. 18, no. 11, pp. 1499–1510, 2008.
- [49] O. Duchenne, I. Laptev, J. Sivic, F. Bach, and J. Ponce, "Automatic annotation of human actions in videos," in *Proc. IEEE Intl. Conf.* on Computer Vision, 2009.
- [50] I. Laptev and P. Pérez, "Retrieving actions in movies," in *Proc. IEEE Intl. Conf. on Computer Vision*, 2007.
- [51] L. Cao, Z. Liu, and T. S. Huang, "Cross-dataset action detection," in *Proc. IEEE Conf. on Computer Vision and Pattern Recognition*, 2010.
- [52] O. Boiman, E. Shechtman, and M. Irani, "In defense of nearestneighbor based image classification," in *Proc. IEEE Conf. on Computer Vision and Pattern Recognition*, 2008.
- [53] C. H. Lampert, "Detecting objects in large image collections and videos by efficient subimage retrieval," in *Proc. IEEE Intl. Conf. on Computer Vision*, 2009.

- [54] P. C. Woodland and D. Povey, "Large scale discriminative training of hidden markov models for speech recognition," *Computer Speech and Language*, vol. 16, no. 1, pp. 25–47, 2002.
- [55] M. Datar, N. Immorlica, P. Indyk, and V. Mirrokni, "Localitysensitive hashing scheme based on p-stable distribution," in *Proc. of Twentieth Annual Symposium on Computational Geometry*, pp. 253– 262, 2004.
- [56] J. Bentley, "Programming pearls," Algorithm Design Techniques, vol. 27, no. 9, pp. 865–871, 1984.
- [57] M. Dikmen and et al., "Surveillance event detection," in Video Evaluation Workshop, 2008.

PLACE PHOTO HERE Ying Wu (SM'06) received the B.S. from Huazhong University of Science and Technology, Wuhan, China, in 1994, the M.S. from Tsinghua University, Beijing, China, in 1997, and the Ph.D. in electrical and computer engineering from the University of Illinois at Urbana-Champaign (UIUC), Urbana, Illinois, in 2001.

From 1997 to 2001, he was a research assistant at the Beckman Institute for Advanced Science and Technology at UIUC. During

summer 1999 and 2000, he was a research intern with Microsoft Research, Redmond, Washington. In 2001, he joined the Department of Electrical and Computer Engineering at Northwestern University, Evanston, Illinois, as an assistant professor. He is currently an associate professor of Electrical Engineering and Computer Science at Northwestern University. His current research interests include computer vision, image and video analysis, pattern recognition, machine learning, multimedia data mining, and human-computer interaction. He serves as associate editors for IEEE Transactions on Image Processing, SPIE Journal of Electronic Imaging, and IAPR Journal of Machine Vision and Applications. He received the Robert T. Chien Award at UIUC in 2001, and the NSF CAREER award in 2003. He is a senior member of the IEEE.

PLACE PHOTO HERE Junsong Yuan (M'08) received the Ph.D. and M.Eng degrees from Northwestern University, Illinois, USA, and National University of Singapore, respectively. Before that, he graduated from the special program for the gifted young in Huazhong University of Science and Technology, Wuhan, P.R.China, and received his B.Eng in communication engineering. During the summer 2008, 2007 and 2006, he was a research intern with the Communication and Collaboration Systems

group, Microsoft Research, Redmond, WA, Kodak Research Laboratories, Rochester, NY, and Motorola Laboratories, Schumburg, IL, respectively. From 2003 to 2004, he was a research scholar at the Institute for Infocomm Research, Singapore.

In Sept. 2009, he joined Nanyang Technological University as a Nanyang assistant professor. His current research interests include computer vision, image and video data mining and content analysis, multimedia search etc. He was a recipient of the Doctoral Spotlight Award from IEEE Conf. on Computer Vision and Pattern Recognition (CVPR'09), a recipient of the Nanyang Assistant Professorship from Nanyang Technological University, and received the Outstanding Ph.D. Thesis award from the EECS department in Northwestern University. He is a member of IEEE and ACM.

PLACE PHOTO HERE Zicheng Liu (SM'05) is a senior researcher at Microsoft Research, Redmond. He has worked on a variety of topics including combinatorial optimization, linked figure animation, and microphone array signal processing. His current research interests include activity recognition, face modeling and animation, and multimedia collaboration. He received a Ph.D. in Computer Science from Princeton University, a M.S. in Operational Research from the Institute of Applied Math-

ematics, Chinese Academy of Science, and a B.S. in Mathematics from Huazhong Normal University, China. Before joining Microsoft Research, he worked at Silicon Graphics as a member of technical staff for two years where he developed a trimmed NURBS tessellator which was shipped in both OpenGL and OpenGL-Optimizer products. He has published over 70 papers in peer-reviewed international journals and conferences, and holds over 40 granted patents. He has served in the technical committees for many international conferences. He was the co-chair of the 2003 ICCV Workshop on Multimedia Technologies in E-Learning and Collaboration, the technical co-chair of 2006 IEEE International Workshop on Multimedia Signal Processing, and the technical co-chair of 2010 International Conference on Multimedia and Expo. He is an associate editor of Machine Vision and Applications journal, and a senior member of IEEE.



Fig. 10. A detection example with performing speed and style variations. The yellow bounding box is the ground truth label of the whole human body action and the red bounding box is our detection of two-hand waving.



Fig. 11. A detection example with large spatial scale changes.



Fig. 12. A detection example with cluttered and moving background, as well as severe partial occlusions.



Fig. 13. A false detection example caused by two individual hand-wavings from two different persons.



Fig. 14. An example of missed detection. Although the action is detected in the first few frames, the whole action instance is still treated as a missed detection due to limited overlap with the ground truth labeling.