

Human Action Recognition Using Average VLBP

Dr. Chidananda H

Department of Computer Science and Engineering,
RYMEC.

Abstract:

This work, we proposed a new method to perceive human actions from visual information. An efficient alternative descriptor is constituted by customizing Volume Local Binary Patterns (VLBP) as Average VLBP. An Average VLBP of static human actions along with its visual information in spatio-temporal domain is determined. A Random Forest classifier is trained to learn histograms of human actions. This work has been tested on the popular and freely accessible KTH and Weizmandataset for human action understanding from input visuals. The outcomes show best instate-of-the-art efficiency in contrast with different strategies.

Keywords —Human Action Recognition, Average Volume Local Binary Patterns, Random Forest, Machine Learning, Action binary patterns, ABP, Human activity recognition.

I. INTRODUCTION

One of the complex and emerging area in computer vision is Human action recognition as this requires characteristics of object, activity and time related data. Moreover, activities have been classified into human activities, interaction between humans, object and human interactions and actions that exist in a group (Aggarwal & Ryoo-2011). Activity recognition is a tough task because of the variations in the background environment like varying view positions or blocked view points and changing backgrounds.

Furthermore, actor may perform the action in a different pattern which may generate lot of variations in the object's actions and its characteristics (Poppe, 2010). This work addresses the issue of perceiving activities executed by a solitary individual, for example boxing, applauding, waving, strolling, running. We propose a basic and effective novel element type, in particular Action Binary Pattern (ABP), which joins static article

appearances just as activity data in the spatio-transient space, in one descriptor. An ABP is processed from three visual information followed by a histogram calculation, prompting an invariance against various video lengths. At long last, the histogram is utilized to become familiar with a Random Forest classifier. The proposed approach is assessed on the single view KTH dataset (Schuldt et al., 2004) and Weizman (Blank et al., 2005; Gorelick et al., 2007) dataset just as on the IXMAS (Weinland et al., 2006; Weinland et al., 2010) dataset for multi-view activity acknowledgment. Few dataset visuals used in this work are shown in fig 1 and fig 2.

Related Work The authors Zhao et al. (Zhao and Pietikainen, 2007) represented dynamic surfaces using principal Local Binary Pattern on Three Orthogonal Planes (LBP-TOP) on the other hand the authors (Mattivi and Shao, 2009; Shao and Mattivi, 2010) used the same LBP-TOP to perceive human actions from visuals. The authors

concluded at 88.19% accuracy on the KTH dataset and by using Extended Gradient LBPTOP with Principal Component Analysis in (Mattivi and Shao, 2009) they achieved at 91.25%. Best outcomes (92.69%) were accomplished by consolidating Extended Gradient LBP-TOP with Dollar's discovery strategy (Shao and Mattivi, 2010). Liu et al. (Liu and Yuen, 2010) defined a boosted Eigen Actions framework which determines a spatiotemporal information saliency map (ISM) by predicting pixel density functions. It has only 81.5% accuracy for the KTH dataset but using the Weizman dataset achieved accuracies up to 98.3%. However, the authors Yeffet et al. (Yeffet and Wolf, 2009), achieved good results with accuracy above 90%. In recent times, Kihl et al. (Kihl et al., 2013) achieved a very good results upto 93.4% by adopting a series of local polynomial approximation of Optical Flow (SoPAF).

Numerous methodologies have disadvantages, for example the measure of highlight lead to ambiguities, can't manage distinctive video lengths or have an enormous component space. To beat these issues, we propose another component that joins the capacities of portraying static article appearances just as activity data, in a one descriptor.

II. METHOD

This Section explores on the two important concepts: Volume Local Binary Patterns and Action Binary Patterns. VLBP have gained popular for portraying some parameters in the spatio-temporal space and are determined from basic LBP.

Our proposed ABP is determined in the X-Y-T space as well and moreover makes the fleeting stride size into account.

In section 2.3, the notable artificial learning approach Random Forest by Leo (Breiman, 2001) have been described.

A. VOLUME LOCAL BINARY PATTERN

A LBP were constituted for first time in (Ojala et al., 1994) for classification of texture. The first LBP is figured in a 3×3 cell by contrasting each dark value with the middle one. In the event that the neighbor values are bigger than the middle one, a 1 is considered to the comparing position, otherwise 0.

An eight bits code-word is determined by processing a 3×3 -LBP. This code-word is deciphered as a binary word and changed over to a decimal number. A histogram of all these decimal numbers is fabricated finally to obtain its pattern.

Since LBP highlights portray only static visuals, they are not appropriate for activity recognition where activity and time data ought to be considered. A Volume Local Binary Pattern (VLBP) is presented in (Zhao and Pietikainen, 2007). It is processed in the spatial and temporal area and able to perceive dynamic surfaces In (Zhao and Pietikainen, 2007), the authors characterize a span around the middle point inside the space-time volume from three nonstop edges to get neighboring pixels as opposed to utilizing a 3×3 cell from one frame. The calculation of a VLBP is like the LBP: if the dark benefit of neighboring voxels inside the space time volume is bigger than that of the voxel's middle, the relating position is appointed to a 1, in any case 0. By figuring a VLBP the code word of bit length 24 would be computed as $2^{24} = 16777216$ number of templates. Like the LBP, a histogram of all happening designs is figured. Frequently, this enormous component pool prompts a few ambiguities. To conquer this issue (Fehr, 2007; Topi et al., 2000) presented a uniform LBP and show that the general measure of LBPs can be decreased to a little subset. Experiments on object recognition show that 90% of every conceivable example have a place with this subset. For our application of activity recognition VLBP were unsatisfactory. Probably, the last component pool contains insufficient data to discover discriminative examples in the spatio-temporal-space.

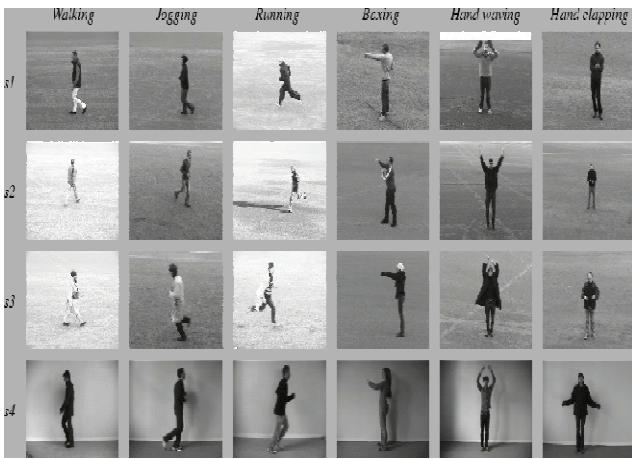


Figure 1: Patterns of the individual perspective of KTH dataset. These images consist of 6 activities performed by different individuals under clustered backgrounds.



Figure 2: Patterns of the individual perspective of Weizman dataset. These images consist of 9 activities performed by different individuals.

Figure 3 depicts VLBP computation. The computed value is 2039583. The 3 successive frames along with computations of its histogram for all the generated patterns are used to compute final VLBP. A random forest classifier is trained by feeding these histogram values.

B. ACTION BINARY PATTERN

This Section explores the attributes and the calculation of our proposed Action Binary Pattern portrayed. Accepting that activity can be recognized by the difference in pixel force esteem, ABPs are

registered from edges of 3 frames and measure the activity between them. Like the Optical Flow (as defined in 1981 by Schunck and Horn,), an ABP depicts qualities of activity. Figure 4 depicts the computation of the above said process.

In this work, frames of one complete half of the activity cycle is considered and every frame is partitioned into sub-volumes and the feature histogram is formed for some interested sub-volume/s by concatenating the sub-volume histograms. Using the sub-volume representation, action and shape are encoded on region-level (sub-volume histogram). To obtain a rough spatial definition of human leg movements, we divide the xyt volume into four regions through the centroid of the silhouette.

This work considered only 50% of one complete cycle of activity. In this , each extracted image frame is further subdivided into its sub-volumes and there final featured histogram is computed by sub-volume histogram concatenation. Utilizing the sub-volume representation, activity and shape are encoded on region-level (sub-volume histogram).

To get a spatial definition of human leg activities, we separate the xyt volume into four sectors through the centroid of the outline. This division generally isolates the hands and legs of the individual. Utilizing more squares would of course permit a more nitty gritty portrayal but would too deliver more neighborhood histograms and make the total histogram longer. The sub-volume division and the arrangement of our featured histogram are outlined in Figure 4. The ABP - are computed from the entire half cycle of the activity and concatenated all the back leg locale sub-volume histograms on each plane.

In the back leg region subvolume, the frames are separated into cells and for three cells at a similar position, the relating estimations of each cell of the starting frame is contrasted and the determi mean estimation of the subsequent casing's inside cell esteem. The inside cell estimation of the subsequent

casing is the normal dim estimation of all the 3X3 cells. This is as appeared in fig 4.

In the event that the gray value inside one cell of the starting frame is bigger than that in the subsequent frame's middle cell, a 1 is considered, else 0. By utilizing a similar strategy, the third frame is contrasted with the subsequent frame.

With respect to calculation of all ABPs in three edges, the number N of inspected designs $C_n(x,y)$, $1 \leq n \leq N$ depends upon the casing size, on the example's size and it's progression size. Each example speaks to the activity between these frames while the binary values, particularly the total number of ones $\|C_n\|$, meant by their entrance astute 1-standard, can be deciphered as the quality of activity. To recognize frail or solid activities, an activity vector $\vec{I} = (i_1, \dots, i_N)^T$ is introduced. \vec{I} is inferred by a quality condition:

$$i_n = \begin{cases} 1, & \|C_n\| \geq S, \\ 0, & \text{Otherwise} \end{cases} \quad (1)$$

Every component of the activity vector is set to 1 if the total number of ones in the example $C_n(x,y)$ is higher than or equivalent to the activity quality limit $S \in \{1 \dots 7\}$. Hence, I contains positions with a specific extent of activity. It depicts the spatial-subordinate changing of intensity values from three frames and in a perfect world just situations with enormous changes in the pixel force are doled out to 1. An ABP descriptor is made by figuring a histogram from all activity vectors of a video. In this way, the histogram is legitimately used to build a component for learning a Random Forest classifier.

In this work, edges of one complete portion of the action cycle is considered and each frame is divided into subvolumes and the component histogram is framed for some intrigued subvolume/s by linking the subvolume histograms. Utilizing the subvolume portrayal, activity and shape are encoded on area

level (subvolume histogram). To acquire a harsh spatial meaning of human leg developments, we isolate the xyt volume into four locales through the centroid of the outline. This division generally isolates the hands and legs of the individual. Utilizing more squares would obviously permit a more itemized depiction yet would likewise create more neighborhood histograms and make the entire histogram longer. The subvolume division and the development of our component histogram are delineated in Figure 4. The ABP highlights are determined from the entire span of an arrangement and connect all the back leg region subvolume histograms on each plane.

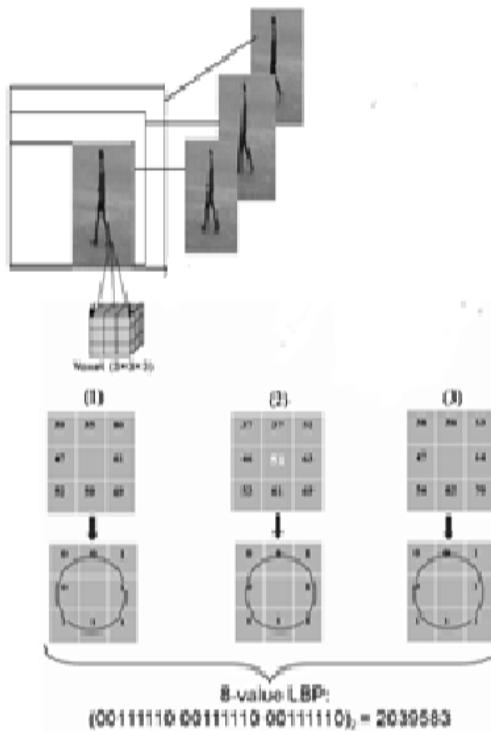


Figure 3: Procedure of computing a Volume Local Binary Pattern.

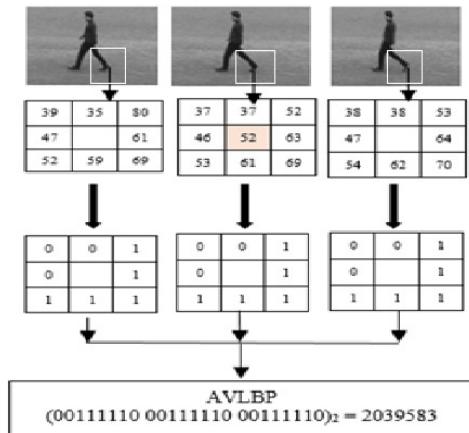


Figure 4: Procedure of computing our proposed ABP in three frames.

Temporal Variations: So as to take in highlights from quick and moderate activities, ABPs are not just figured from three consistent casings. Clearly just quick activities could be perceived by analysing parameters of frame sequences. Furthermore, four spatial scale steps are defined and ABPs are registered by joining these means for moving the pattern through the space-time volume. A period venture of $ts = 1, 2, 3, 4$ was exactly picked. For the instance of $ts = 1$, an ABP is registered from three consistent frames. Consistently outline is gathered for $ts=2$. Separately, for $ts = 3, 4$ every third or fourth frame was picked. Rather than making a solitary histogram that can portray quick activities, four histograms are made to describe diverse sort of activities. These histograms are connected and used to gain proficiency with a Random Forest classifier.

C. RANDOM FOREST

In this section, a short clarification of the hypothesis behind Random Forest is given. Arbitrary Forests were created by Leo (Breiman, 2001) and join packing (Breiman, 1996) with an irregular element choice proposed by (Ho, 1995; Ho, 1998) and Amit (AmitandGeman, 1997). A Random Forest comprises of an assortment of CART-like choice trees ht , $1 \leq t \leq T$:

$$\{h(\vec{x}, \theta_t)_{t=1, \dots, T}\}$$

Where $\{\theta_t\}$ is a bootstrap test from the training data. Each tree makes a choice on a class for the information \vec{x} : The class probabilities are assessed by greater part casting a ballot and used to compute the example's name $y(\vec{x})$ as for a given component vector \vec{x}

$$y(\vec{x}) = \underset{c}{\operatorname{argmax}} \left(\frac{1}{T} \sum_{t=1}^T F_{h_t(\vec{x}) = c} \right) \quad (2)$$

The decision function $h_t(\vec{x})$ returns the result class c of one tree with the indicator function F :

$$F_{h_t(\vec{x}) = c} = \begin{cases} 1, & h_t(\vec{x}) = C, \\ 0, & \text{Otherwise} \end{cases} \quad (3)$$

Random Forest has an high classification accuracy and can manage enormous informational indexes for different classes with extraordinary time efficiency.

D. CLASSIFICATION

Input descriptors are classified by passing them down each tree until a leaf hub is reached. The outcome class is defined by each leaf hub and the final choice is controlled by taking the class having the most votes (dominant part vote), see Equation (2).

III. EXPERIMENTAL RESULTS

KTH: AVLBP and ABPs are assessed on the notable and openly accessible KTH dataset (Schuldt et al., 2004) comprising of six classes of activities. Each activity is performed by 25 people in four distinct situations. The KTH dataset comprises of 599 recordings. Like (O'Hara and Draper, 2012), a fixed position jumping box with a transient window of 24 frames is chosen, in light of explanations by (Lui et al., 2010). Probably, fewer frames is sufficient (Schindler and Van Gool, 2008). Yet, in this work, just half complete pattern of movement is

sufficient as the edges will rehash in the staying half cycle. Moreover, the first preparing/testing parts from (Schuldt et al., 2004) are utilized. A few methodologies for processing VLBP values were tried. Two unique neighborhoods (eight and four qualities) were analyzed, the influence of various histogram goes just as the distinction between frame-by-frame learning and multi-frame learning has been computed. Best outcomes were accomplished by processing a 4-esteem VLBP with multi-outline learning (one histogram for all edges of a video is made) and a histogram scope of 400 canisters. Figure 5(a) shows the disarray framework with a normal precision of 92.83% for the KTH dataset.

A. EVALUATION FOR ACTION BINARY PATTERNS

For determining an ABP, the activity strength threshold must be balanced. This feature is very important that influences the accuracy of the ABP. Besides, ABPs are more sensitive to various picture sizes. In this Section, we tried the influence of these boundaries and contrast the outcomes with a few late methodologies.

| | Box | Walk | Run | Jog | Wave | Clap |
|------|------|------|------|------|------|------|
| Box | 0.95 | 0 | 0 | 0 | 0.02 | 0 |
| Walk | 0 | 1 | 0 | 0 | 0 | 0 |
| Run | 0 | 0 | 1 | 0 | 0 | 0 |
| Jog | 0 | 0.02 | 0.02 | 0.96 | 0 | 0 |
| Wave | 0.12 | 0 | 0 | 0 | 0.82 | 0.12 |
| Clap | 0.06 | 0 | 0 | 0 | 0.13 | 0.84 |

(a)

| | Box | Walk | Run | Jog | Wave | Clap |
|------|------|------|------|------|------|------|
| Box | 1 | 0 | 0 | 0 | 0.02 | 0 |
| Walk | 0 | 0.99 | 0 | 0.03 | 0 | 0 |
| Run | 0.03 | 0.03 | 0.93 | 0.03 | 0.02 | 0.05 |

| | | | | | | |
|------|---|------|---|------|------|------|
| Jog | 0 | 0.02 | 0 | 0.96 | 0 | 0 |
| Wave | 0 | 0 | 0 | 0 | 0.83 | 0.12 |
| Clap | 0 | 0 | 0 | 0 | 0.13 | 0.89 |

(b)

Figure 5: Achieved accuracy: (a) 89.81% using VLBP, (b) 91.83% using ABP.

B. INFLUENCE OF THE ACTION STRENGTH THRESHOLD

A brief clarification about the activity strength threshold $S \in \{1...7\}$ is explained in section 2.2. Table 1 shows the acknowledgment precision when the limit fluctuates for a frame size of 75×150 . At the point when the limit is bigger than seven, there will be less non-zero values in the ABP histogram. As recorded in table1, perceiving exactness increases as the limit increases. The highest exactness 91.32% is accomplished by utilizing a limit of five, prompting the assumption that this worth is entirely appropriate for the task of human activity acknowledgment.

Table1: Action thresholds and Average accuracies for ABP. Better accuracy achieved with a threshold of 5.

| Threshold | Average accuracy (%) |
|-----------|----------------------|
| 1 | 75.60 |
| 2 | 81.62 |
| 3 | 86.40 |
| 4 | 91.21 |
| 5 | 89.23 |
| 6 | 92.24 |
| 7 | 90.32 |

In contrast to the VLBP, ABP doesn't do well in perceiving running. ABP is less efficient to recognise quick activities compared to VLBP as quick activities sometimes produce different descriptors which causes difficult to recognise the activity. Be that as it may, ABPs are delicate to frail activities. Overall results from test shows inaccurate recognition generally occur on some activities like strolling, running, boxing on running, waving, applauding, as appeared in figure 5(b).

Table 2: Threshold values with Average accuracy for ABP. A threshold of 5 is leading to the best accuracy.

| Threshold | Average accuracy (%) |
|-----------|----------------------|
| 1 | 74.33 |
| 2 | 80.59 |
| 3 | 85.38 |
| 4 | 90.33 |
| 5 | 91.32 |
| 6 | 91.24 |
| 7 | 90.32 |

Table 3: Comparison of other methods with our method using the KTH dataset.

| Threshold | Average accuracy (%) |
|-------------------------|----------------------|
| (Kihl et al., 2013) | 91.5 |
| (Yeffet and Wolf, 2009) | 90.1 |
| (Laptev et al., 2008) | 91.8 |
| Our Method (ABP) | 92.11 |

C. COMPARISON TO STATE-OF-THE-ART METHODS

In this Section, we contrast the proposed technique with a few state-of-the-art works and show that the ABPs are a very efficient descriptor for activity acknowledgment. KTH: The ABP accomplishes an accuracy of 92.11% on the KTH dataset. Table 4 reports the accuracy of our proposed ABP in contrast with different strategies. Activity Binary Patterns arrive at most noteworthy precision for unique training-/testing split and is just marginally lower than the best outcome with cross-approval. Confusion matrix is shown in figure 5.

IV. CONCLUSIONS AND FUTURE WORK

In this paper a novel component type, to be specific Action Binary Patterns (ABP) are proposed. ABPs consolidate the benefits of Volume Local Binary Patterns to produce an Average VLBP to accumulate static object information and Optical Flow to get activity data. An ABP is processed from three frames with a temporal shifted sliding window. The subsequent histograms are utilized to become familiar with a Random Forest classifier. The proposed features are assessed on the notable, freely accessible KTH dataset. The outcomes

exhibit cutting edge accuracies in contrast with Volume Local Binary Patterns.

Future Work Our arrangements for future work are to assess VLBP on more complex datasets Furthermore, we proposed to arrange of the manual alter of the action constrain by showing an entropy work that extracts designs with more discriminative characteristics. On the other hand, we propose to encode more information into the case.

For instance, all the transient action patterns may be joined into one final histogram and this requires an extra care in selecting the perfect cell size if we consider an Average VLBP. In this paper we proposed to prepare an Average VLBP in a 3×3 cell however the results may be made strides by taking other cell sizes like 5×5 or 7×7 .

REFERENCES

- Aggarwal, J. and Ryoo, M. (2011). Human activity analysis: A review. *ACM Computing Surveys*, 43(3):16:1– 16:43.
- Amit, Y. and Geman, D. (1997). Shape quantization and recognition with randomized trees. *Neural computation*, 9(7):1545–1588.
- Blank, M., Gorelick, L., Shechtman, E., Irani, M., and Basri, R. (2005). Actions as space-time shapes. In Computer Vision (ICCV), 10th International Conference on, pages 1395–1402.
- Breiman, L. (1996). Bagging predictors. In Machine Learning, volume 24, pages 123–140.
- Breiman, L. (2001). Random forests. *Machine learning*, 45(1):5–32.
- Fehr, J. (2007). Rotational invariant uniform local binary patterns for full 3d volume texture analysis. In Finnish signal processing symposium (FINSIG).
- Gorelick, L., Blank, M., Shechtman, E., Irani, M., and Basri, R. (2007). Actions as space-time shapes. *Pattern Analysis and Machine Intelligence (PAMI)*, IEEE Transactions on, 29(12):2247–2253.
- Ho, T. K. (1995). Random decision forests. In Document Analysis and Recognition, 1995., Proceedings of the Third International Conference on, volume 1, pages 278–282. IEEE.
- Ho, T. K. (1998). The random subspace method for constructing decision forests. *Pattern Analysis and Machine Intelligence, IEEE Transactions on*, 20(8):832– 844.
- Horn, B. K. and Schunck, B. G. (1981). Determining optical flow. *Artificial Intelligence*, 17.
- Jhuang, H., Serre, T., Wolf, L., and Poggio, T. (2007). A biologically inspired system for action recognition. In Computer Vision (ICCV), 11th International Conference on, pages 1–8. IEEE.
- Kihl, O., Picard, D., Gosselin, P.-H., et al. (2013). Local polynomial space-time descriptors for actions classification. In International Conference on Machine Vision Applications.
- Laptev, I., Marszalek, M., Schmid, C., and Rozenfeld, B. (2008). Learning realistic human actions from movies. In Computer Vision and Pattern Recognition, (CVPR). IEEE Conference on.
- Li, R. and Zickler, T. (2012). Discriminative virtual views for cross-view action recognition. In Computer Vision and Pattern Recognition, (CVPR). IEEE Conference on.
- Li, W., Yu, Q., Sawhney, H., and Vasconcelos, N. (2013). Recognizing activities via bag of words for attribute dynamics. In Computer Vision and Pattern Recognition, (CVPR). IEEE Conference on, pages 2587–2594.

16. Lin, Z., Jiang, Z., and Davis, L. S. (2009). Recognizing actions by shape-action prototype trees. In Computer Vision (ICCV), 12th International Conference on, pages 444–451. IEEE.
17. Liu, C. and Yuen, P. C. (2010). Human action recognition using boosted Eigen actions. Image and vision computing, 28(5):825–835.
18. Liu, J., Luo, J., and Shah, M. (2009). Recognizing realistic actions from videos “in the wild”. In *Computer Vision and Pattern Recognition, 2009. CVPR 2009. IEEE Conference on*, pages 1996–2003. IEEE.
19. Lui, Y. M., Beveridge, J., and Kirby, M. (2010). Action classification on product manifolds. In Computer Vision and Pattern Recognition, (CVPR). IEEE Conference on.
20. Marszałek, M., Laptev, I., and Schmid, C. (2009). Actions in context. In Computer Vision and Pattern Recognition, (CVPR). IEEE Conference on.
21. Mattivi, R. and Shao, L. (2009). Human action recognition using lbp-top as sparse spatio-temporal feature descriptor. In Computer Analysis of Images and Patterns (CAIP).
22. O’Hara, S. and Draper, B. (2012). Scalable action recognition with a subspace forest. In Computer Vision and Pattern Recognition, (CVPR). IEEE Conference on.
23. Ojala, T., Pietikainen, M., and Harwood, D. (1994). Performance evaluation of texture measures with classification based on kullback discrimination of distributions. In Pattern Recognition. Proceedings of the 12th IAPR International Conference on.
24. Poppe, R. (2010). A survey on vision-based human action recognition. Image and Vision Computing, 28(6):976 – 990.
25. Schindler, K. and Van Gool, L. (2008). Action snippets: How many frames does human action recognition require? In Computer Vision and Pattern Recognition, (CVPR). IEEE Conference on.
26. Schuldt, C., Laptev, I., and Caputo, B. (2004). Recognizing human actions: a local svm approach. In Pattern Recognition. (ICPR). Proceedings of the 17th International Conference on.
27. Shao, L. and Mattivi, R. (2010). Feature detector and descriptor evaluation in human action recognition. In Proceedings of the ACM International Conference on Image and Video Retrieval.
28. Tian, Y., Sukthankar, R., and Shah, M. (2013). Spatiotemporal deformable part models for action detection. In Computer Vision and Pattern Recognition (CVPR). IEEE Conference on.
29. Topi, M., Timo, O., Matti, P., and Maricor, S. (2000). Robust texture classification by subsets of local binary patterns. In Pattern Recognition. (ICPR). Proceedings of the 15th International Conference on.
30. Wang, Z., Wang, J., Xiao, J., Lin, K.-H., and Huang, T. (2012). Substructure and boundary modeling for continuous action recognition. In Computer Vision and Pattern Recognition, (CVPR). IEEE Conference on.
31. Weinland, D., Ozysal, M., and Fua, P. (2010). Making action recognition robust to occlusions and viewpoint changes,. In European Conference on Computer Vision (ECCV).