

Singapore Management University

Institutional Knowledge at Singapore Management University

Research Collection School Of Computing and Information Systems

School of Computing and Information Systems

12-2014

Human action classification based on sequential bag-of-words model

Hong LIU

Qiaoduo ZHANG

Qianru SUN

Singapore Management University, qianrusun@smu.edu.sg

Follow this and additional works at: https://ink.library.smu.edu.sg/sis_research



Part of the [Computer Engineering Commons](#), and the [Software Engineering Commons](#)

Citation

LIU, Hong; ZHANG, Qiaoduo; and SUN, Qianru. Human action classification based on sequential bag-of-words model. (2014). *Proceedings of the 2014 IEEE International Conference on Robotics and Biomimetics (ROBIO 2014), Bali, December 5-10*. 2280-2285.

Available at: https://ink.library.smu.edu.sg/sis_research/4461

This Conference Proceeding Article is brought to you for free and open access by the School of Computing and Information Systems at Institutional Knowledge at Singapore Management University. It has been accepted for inclusion in Research Collection School Of Computing and Information Systems by an authorized administrator of Institutional Knowledge at Singapore Management University. For more information, please email cherylds@smu.edu.sg.

Human Action Classification Based on Sequential Bag-of-Words Model

Hong Liu, Qiaoduo Zhang and Qianru Sun [†]

Abstract—Recently, approaches utilizing spatial-temporal features have achieved great success in human action classification. However, they typically rely on bag-of-words (BoWs) model, and ignore the spatial and temporal structure information of visual words, bringing ambiguities among similar actions. In this paper, we present a novel approach called sequential BoWs for efficient human action classification. It captures temporal sequential structure by segmenting the entire action into sub-actions. Each sub-action has a tiny movement within a narrow range of action. Then the sequential BoWs are created, in which each sub-action is assigned with a certain weight and salience to highlight the distinguishing sections. It is noted that the weight and salience are figured out in advance according to the sub-action’s discrimination evaluated by training data. Finally, those sub-actions are used for classification respectively, and voting for united result. Experiments are conducted on UT-interaction dataset and Rochester dataset. The results show its higher robustness and accuracy over most state-of-the-art classification approaches.

I. INTRODUCTION

Automatic human action classification is of significant use in many applications, such as intelligent surveillance, content-based video retrieval and human-computer interaction. It has been researched for years, but remains a very challenging task. One of the most difficult problems is to distinguish between actions with high inter-ambiguities. Regarding an action as a connection of sub-actions, some action classes consist of similar sub-actions, which greatly increase the difficulty of classification (Fig. 1).

Recently, spatial-temporal feature based approaches (e.g. [1, 2, 3]) have been widely used in human action analysis and achieved promising results. The local features are regarded as visual words, then each action is described as a single descriptor using bag-of-words (BoWs) model. Although BoWs model is popular, it has an essential drawback of only focusing on the number of words but ignoring the spatial-temporal information. This results in ambiguities between

This work is supported by National Natural Science Foundation of China (NSFC, No.60875050, 60675025, 61340046), National High Technology Research and Development Program of China (863 Program, No.2006AA04Z247), Science and Technology Innovation Commission of Shenzhen Municipality (No.201005280682A, No.JCYJ20120614152234873), Specialized Research Fund for the Doctoral Program of Higher Education (No. 20130001110011).

Hong Liu is with the Engineering Lab on Intelligent Perception for Internet of Things(ELIP), and the Key Laboratory of Machine Perception, Peking University, Beijing, 100087 CHINA. hongliu@pku.edu.cn

Qiaoduo Zhang is with the Engineering Lab on Intelligent Perception for Internet of Things(ELIP), Shenzhen Graduate School, Peking University, Shenzhen, 518055 CHINA. qiaodtg@sz.pku.edu.cn

Qianru Sun is with the Engineering Lab on Intelligent Perception for Internet of Things(ELIP), Shenzhen Graduate School, Peking University, Shenzhen, 518055 CHINA. qianrusun@sz.pku.edu.cn



Fig. 1. Examples of two confusing actions: “punching” and “kicking”. Segmenting them into 5 sections, the former sections are nearly the same. This leads to large ambiguity between their BoWs histograms. But if we focus on the salient parts of the actions, the classification would be much easier.

different classes of actions. For example, some similar sub-actions happen at different relative period of time, such as stand up at the beginning of an action or at the end. BoWs model is incapable to distinguish them. Moreover, BoWs model handles all the visual words equally and can not lay stress on the most distinguish parts. The higher proportion of similar sub-actions between classes, the more difficult to do the classification using original BoWs. Hence, considering temporal series of BoWs model is of great importance and necessity.

However, only a small part of previous works focus on capturing words’ temporal relationships. In some approaches [1, 2, 3], spatial-temporal correlations of local features were learned as neighborhoods or correlograms to capture the words’ spatial-temporal relationships. While these approaches are still too local to capture long-term relationship between words. Other works [4, 5, 6] counted the co-occurrence between words, but they limited the scope within a small period of time. Recently Ryoo [7] represented an activity as an integral histogram of spatial-temporal features, efficiently modeling how feature distributions change over time, but it could not deal with ambiguities between classes those have similar sub-actions at same relative time. Glaser et al. [8] incorporated temporal context in BoWs models to capture the contextual sequence between words but it still could not focus on the distinctive parts of the actions.

Instead of directly including time information in visual words, we take account of time dimension by segmenting the entire action into small sections. Each section is called a sub-action. The i th sub-actions of all actions compose the i th sub-section. Then sequential BoWs model is used for video representation, the classification is described as a series of sub-classes classification problems. In this way, we can apply our approach to original BoWs and spatial-

improved BoWs, such as approaches in our previous works [5, 6]. Satkin et al. [13] extracted the most discriminative portion of a video for training according to the accuracy of a trained classifier. Inspired by this, the discriminative parts of the action is emphasised by assigning them high weights and salience values. The weight indicates the probability of a sub-action belonging to a certain class and the salience means how much a sub-section is distinguished from others. Then the effect of similar sub-actions is minimized, so we can focus on the difference between classes (salient parts in Fig. 1). On the one hand, if similar sub-actions happen at different relative time, they will be in different sub-section and classified separately. On the other hand, if they can not be separated, then low salience values are given to them to reduce their influence. Finally, sub-actions are classified separately and the results will be gathered together through voting. The experiments are implemented on UT-interaction dataset and a more challenging Rochester dataset. The results show our approach can achieve robust and accuracy beyond most related approaches.

The rest of the paper is organized as follows. In Section II, we introduce the reason for using sub-actions and illustrate the framework of our approach. Section III and Section IV describe the segmentation and classification approach respectively. In Section V, we conduct experiments on UT-Interaction dataset and Rochester dataset and compare our approach with other BoWs based approaches. Finally, conclusions are drawn in Section VI.

II. FRAMEWORKS

Human movement can be described at various levels of complexities [15]. Usually, an activity refers a whole expression of movement such as “play tennis”. An action is the element of activity such as “running” or “jumping”. It is often short and represents less motion information. Finally, an action primitive is a very short period of action that can not be performed individually. Action primitives are components of actions and actions are as well components of activities.

In general, action classification is performed at the first two levels. Many different action classes are unavoidable to share similar or same action primitives. This largely increases the difficulty to distinguish different classes. However, performing action classification at primitive level is also impractical. The reason is that there are innumerable action primitives due to the flexibility of human movement. Moreover, different action classes may be composed by different numbers of action primitives, it is not suitable to perform uniform preprocessing for all action classes. To solve this problem, we define sub-action instead of action primitive. Sub-action is also a small period of action, but it is not previously defined or fixed for each action classes. Sub-action is segmented according to a certain action automatically. All the action classes to be classified will have same number of sub-actions. This approach can not only solve the problem mentioned above but also be faster, more flexible and adaptive.

As shown in Fig. 2, first local features are extracted from action videos. The videos are then treated as volumes of visual words. To extract sub-actions, we chop the volumes into small clips according to the intensity of actions. Distances between clips are accumulated. Then segment the accumulated distance into equal parts, i.e. the sub-actions. Sequential histograms are created for each action respectively to describe the action. Before classification, similarity between sub-actions in the same sub-section is figured out with training data. These similarities can be regarded as weights for voting. Meanwhile salience value for each section is figured out according to the pre-classification accuracy. After sub-section classification, a vote scheme is conducted finally. Category of a test instance A is decided by the equation below,

$$A = \max_{i \in C} \sum_{j=1}^{N_s} \omega_j(i, c_j) s_j(c_j), \quad (1)$$

where C denotes all the possible categories, N_s represents the number of sub-sections, c_j denotes A 's sub-classification result in sub-section j and $\omega_j(i, c_j)$ is the weight of the j th sub-action belonging to class i when it is classified to subclass c_j . The last $s_j(c_j)$ stands for the saliency of the j th sub-action. The instance will be classified to the category with the highest score.

III. ACTION SEGMENTATION

Our approach takes advantage of local spatio-temporal features to represent actions. After extracting local features into visual words from video sequence, we conduct a two-stage segmentation to avoid acquiring inequality sub-actions caused by the scale, range, rate or other individual difference between actors. Optimal segmentation could narrow the classification and disambiguate similar sub-actions occur at different relative time. We briefly introduce the visual words extraction approach in III-A, and in III-B, our segmentation approach to acquire sub-actions is presented.

A. Visual Words Extraction

Local feature based representation is widely used in human action analysis for its resistant to clutters and noise. Firstly, spatial-temporal interest points (STIPs) are detected, small cuboid is extracted around each interest point. Then descriptors are generated to represent local information. There are many different local detectors and descriptors being proposed. Since the average length of sub-action is short, our approach avoid using those detecting approaches which are too sparse. Dense detector [10, 11] are all good choices.

Then k-means clustering algorithm is used to build a visual word dictionary and each feature descriptor is assigned to the closest word in the vocabulary. Finally, a video with N frames is described as a sequence of frames with visual words:

$$video = [f_1, f_2, \dots, f_N], \quad (2)$$

where,

$$f_i = [w_{1_i}, w_{2_i}, \dots, w_{n_i}], \quad (3)$$

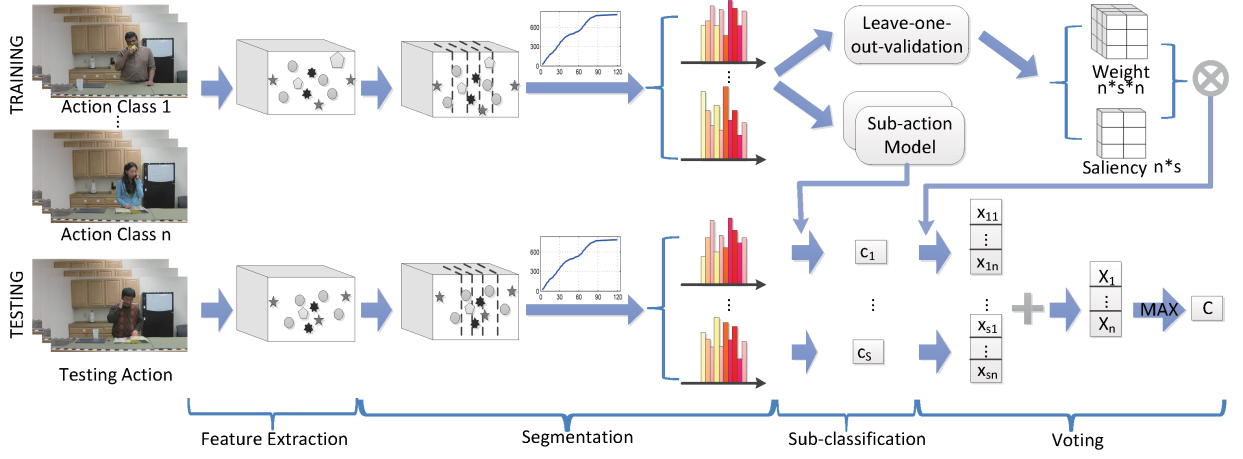


Fig. 2. Framework of our approach. Firstly, videos are transformed into volumes of visual words. Secondly, a two-stage segmentation is conducted. Stage-1 chop volumes into clips according to the density of words, stage-2 divide accumulated distance between clips into equal parts. Thirdly, normal classification is done to each sub-section respectively and results are recorded as c_j . Finally, the final result is decided by a voting process.

n_i is the word number in the i th frame.

B. Sub-action Segmentation

To segment actions efficiently, we should achieve two goals. First, ensure the sub-actions in the same sub-section of the same action class are of the same type, ignoring the speed differences between actors. Second, all sub-actions should capture enough motion information for classification. The two-stage segmentation can reach the above goals nicely, which is detailed below.

1) *Chop Clips*: In first segmentation stage, the entire video is chopped into normalized clips with approximately equal numbers of feature points. The k th clip ends up in the x_k th frame, x_k should satisfy the equation below:

$$\sum_{i=1}^{x_k-1} n_i < \left(\sum_{i=1}^N n_i \right) * k / N_c \leq \sum_{i=1}^{x_k} n_i, \quad (4)$$

where N_c is the clip number.

Here we use point density to chop clips for further segmentation instead of using number of frames to ensure there is enough motion information in all clips. Generally, dense feature points infer strenuous action. The intensity of the action may be changed over the entire action process. Some of the action primitives maybe moderate, so there could be few interesting points in these frames and result in insufficient information. Moreover, it also balances the speed difference between different actors or actions. We compute a histogram of spatial-temporal word occurrence for each video clip, the k th clip is described as:

$$h_k = \text{hist}([w | w \in f_j, x_{k-1} < j < x_k]). \quad (5)$$

2) *Segment Sub-actions*: In second segmentation stage, motion range is calculated to segment sub-actions. The motion range is measured with χ^2 distance between neighbor clips. The distance from clip i to clip j is:

$$\text{dist}(i, j) = \sum_{k=i}^{j-1} \chi^2(h_k, h_{k+1}). \quad (6)$$

Then accumulated distance of the whole action is divided into equal parts and the motion range of sub-actions is figured out:

$$T = \text{dist}(1, N_c) / N_s. \quad (7)$$

It is used as the threshold to segment clip series. In fact, the distance between clips infers not only the range of motion, but also the changing extent of the action primitives' type. Finally, the sub-actions are segmented by the equation below:

$$\text{dist}(i, j) \leq T < \text{dist}(i, j + 1), \quad (8)$$

where i, j are the beginning and ending clips of the sub-section. A stable segmentation over classes is achieved by concatenating the adjacent clips together.

A segmentation example is illustrated in Fig. 3, and the corresponding accumulated distance curves are shown in Fig. 4. By segmenting action like this, we can eliminate some of the speed and range differences between instances. Therefore to a certain class, same sub-actions could be basically segmented to the same sub-section ignoring the delicate length difference between instance. Sequential BoWs is formed for each sub-action respectively for weight calculation and sequential classification. All segmentation steps are shown in Algorithm 1. The algorithm focuses on sub-action segmentation to achieve equally distribution among different sections.

IV. SUB-CLASSIFICATION AND VOTING

A pre-classification using training data is conducted to figure out the weight and salience value for each sub-action. The weight shows the sub-action's discrimination with other sub-classes, while the salience indicates the sub-action's importance within it's own class.

A. Weight Calculation

Similar sub-actions may occur in the same sub-section of different classes, hence the result of directly voting would be poor. The pre-classification result $M(i, k)$ represents the percentage of sub-class i be classified to sub-class k . If the j th sub-action is classified as sub-class k , the probability for

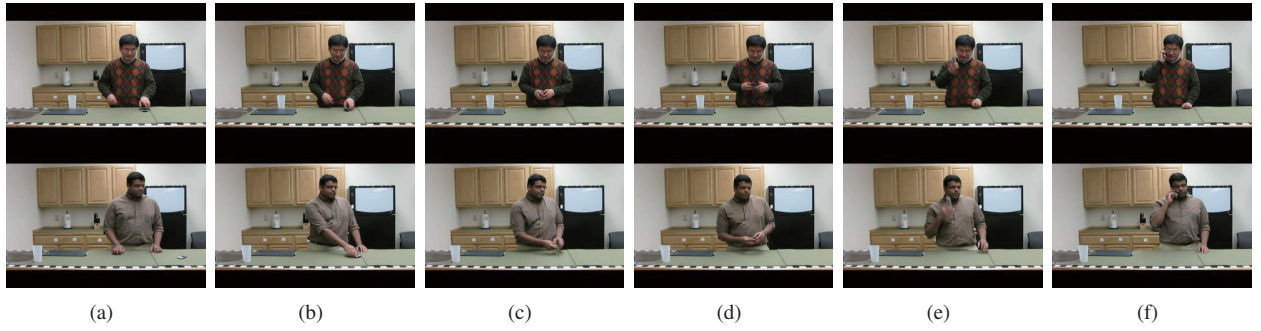


Fig. 3. The key images of “answer phone” acted by two actors. (a, b, c, d, e, f) correspond to six segmenting points A, B, C, D, E, F in Fig. 4. (a, f) are the start and end of the video. In the first section, actions change from (a) to (b), both actors’ hands go forward to get the phone. In the second section (b) to (c), the phones are taken closer. In the third section (c) to (d), actors open the clamshell phones. In the fourth section (d) to (e), actors raise the phones. In the final section (e) to (f), actors listen to the phones. Both final sections are longer because actors move little.

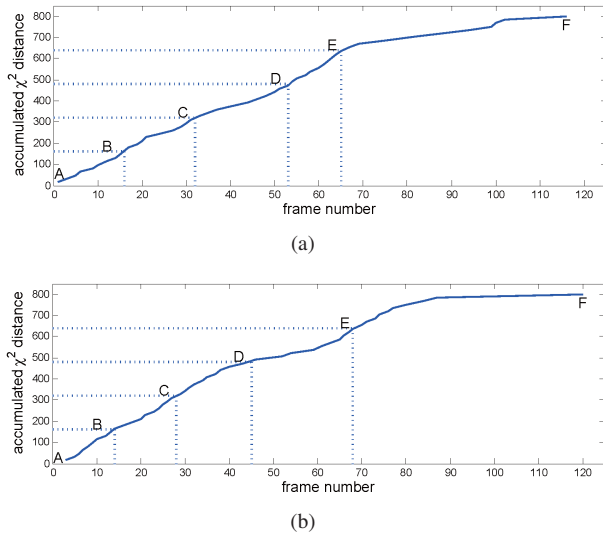


Fig. 4. Example of accumulated χ^2 distance changing over time. (a) represents the upper actor in Fig. 3 and (b) represents the lower one. The accumulated distance curves are similar although they are acted by different actor and rate. B, C, D, E are four equal segmenting points. Their corresponding images are shown in Fig. 3. At each point, the actions are almost executed to the same ratio.

this sub-action belonging to class i is calculated below:

$$\omega_j(i, k) = P(i|c_j = k) = \frac{M_j(i, k)}{\sum_{l \in C} M_j(l, k)} \quad (9)$$

This value can be regarded as weight to eliminate the ambiguity between sub-classes. The difference between weights shows the dissimilarity between sub-classes. The smaller the difference is, the higher the similarity shows. Similar sub-actions give approximately equally increase to their own category when voting, so the classification result is barely changed.

B. Saliency Calculation

The saliency value for each sub-action is figured out to differ the importance of each sub-action within the class. In some sub-section, sub-actions are at low similarity so the classification accuracy would be high. We assign high saliency scores to sub-actions in such sub-sections. While for some other sub-sections with large ambiguities, classification

Algorithm 1 Sub-action Segmentation

Require: Visual word set W , their location $sub = x, y, t$, number of clips N_c , number of sub-actions N_s

Ensure: Sub-action histograms H

- 1: Sort W by the order of t in sub
- 2: $b = \text{sizeof}(W)/N_c$;
- 3: **for** $i = 1$ to N_c **do**
- 4: $h(i) = \text{hist}(W((i-1)*b, i*b))$;
- 5: **if** $i \neq 1$ **then**
- 6: $X(i-1) = \chi^2(h(i-1), h(i))$;
- 7: **end if**
- 8: **end for**
- 9: $T = \text{sum}(X)/N_s$; $temps = 0$; $is = 1$; $seg(1) = 1$;
- 10: **for** $i = 1$ to N_c **do**
- 11: $temps = temps + X(i)$;
- 12: **if** $temps > T * is$ **then**
- 13: $seg(is+1) = i$; $is = is + 1$;
- 14: **end if**
- 15: **end for**
- 16: **for** $i = 1$ to N_s **do**
- 17: $H(i) = \text{sum}(h(seg(i)) : h(seg(i+1) - 1))$
- 18: **end for**

maybe hard. To reduce the effect of these sub-section, low saliency scores are assigned. For a single category, we can figure out its classification accuracy in each sub-section via training sets. But at the testing time, we can not ensure the category for an action, so an average value is used as below:

$$s_j(k) = \sum_{l \in C} \frac{M_j(l, l)}{\sum_{i=1}^{N_s} M_i(l, l)} * \omega_j(l, k) \quad (10)$$

Fig. 5 shows an example of calculated saliency maps for UT-Interaction scene-1, scene-2 and the difference between them. To each column, the saliency value is proportional to the sub-action’s distinctiveness. For example the first column (shake hands), the shaking parts of the action are more salient. But to the fourth column (point), the hole process is quite different from other actions, so its saliency is evenly distributed. Although the environment of the two scenes is different, their saliency maps are quite similar. The saliency

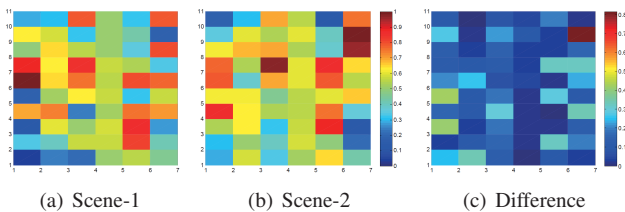


Fig. 5. Saliency maps calculated on UT-Interaction scene-1, scene-2 and the difference between them. Each grid represents a sub-action’s saliency. The horizontal axis represents different classes, from left to right are “shake hands”, “hug”, “kick”, “point”, “punch”, and “push”. The vertical axis represents sub-actions align by time.

ranges from 0 to 1, and the average difference between them is 0.1777. Except the red grid on the right top of (c) caused by overact in scene-2, most difference is less than 0.2. This means our segmentation and saliency calculation approach is robust to the inner-class variation and environmental change.

Finally, category of a test instance is decided by the multiplication of those two scores (9), (10) calculated above, as in (1). By combining these two weights, the effect of ambiguous sub-actions in the same sub-section is weakened, the salient points at temporal dimension are stressed. We can focus the classification on the distinguishing action parts rather than being confused by those similar parts.

V. EXPERIMENTS

In this section, the proposed classification approach is implemented and evaluated on two challenging datasets, UT-Interaction [12] and Rochester [9]. Fig. 6 shows some actions from these two datasets. We compare our results to state-of-the-art classification approaches and confirm the advantages of our approach to distinguish similar actions.

The segmented version of the UT-Interaction dataset contains videos of six types of human actions [16]. All besides “pointing” are interactive activities and are performed by two actors. There are two sets in the dataset performed in different environment (Fig. 6 (a) (b)). One is relatively simple, which is taken in a parking lot. The other is more complex with moving trees in the background. Each activity is repeated for 10 times per set by different actors, so there are totally 120 videos. The videos are taken with camera jitters and/or passerby and have been tested by several state-of-the-art approaches [6, 7]. Rochester dataset contains 150 videos of 10 types of actions. Each category is performed by five actors, repeated for three times in the same scenario. The inter-class ambiguities of these two datasets are both large.

We use the cuboid feature detector [10] as local STIPs detector for its simplicity, fastness and generality [14]. A 100-dimensional Dollar’s gradient descriptor [10] and a 640-dimensional 3D-SIFT descriptor [17] are used respectively to describe the STIP-centered cuboids. Other feature detectors and descriptors can also be used to acquire features. For both datasets, the cuboid size is $w = h = 1 \text{ pixels}$, $\tau = 2 \text{ frames}$, threshold is 0.0002 when using gradient descriptor. When using 3D-SIFT, $w = h = 2 \text{ pixels}$, $\tau = 3 \text{ frames}$, the threshold is 0.0001. After extracting features from videos,



Fig. 6. Examples in datasets

k-means clustering is used to transform them into visual words. To UT-Interaction dataset, the cluster numbers in two sceneries are 90 and 140. The cluster number in Rochester is 900. Classification is conducted using 1-NN since the performance of SVM and 1-NN is quite close [14] while 1-NN is faster.

During the segmentation, we use χ^2 distance to measure motion range. The clip number N_c and the sub-action number N_s are decided by iteration tests. N_c is a bit smaller than the frame number of the shortest video, N_s is associated to the action’s complexity. We set $N_c = 140$, $N_s = 90$ for UT-Interaction scene-1, $N_c = 120$, $N_s = 90$ for UT-Interaction scene-2, and $N_c = 150$, $N_s = 20$ for Rochester. Although the average video length of Rochester is longer than that of UT-Interaction, the actions are slow and some moments are even motionless, so the sub-action number of Rochester is smaller. Pre-classification is done in the training set to figure out the weight and saliency. The leave-one-sequence-out cross validation setting is used in all experiments. All confusion matrices are the average of 10 runs since k-means clustering is randomly initialized.

Confusion matrices of UT-Interaction scene-1 and Rochester are shown in Fig. 7. In each column, our sequential BoWs approach using different descriptors is compared with Dollar’s original BoWs. The result of our approach is much better than original BoWs. On UT-Interaction scene-1, errors among “kick”, “punch” and “push” are most obvious in (a). These three actions share a lot of same sub-actions. Moreover their unique sub-actions (stick out one’s arms/leg to conduct a hit/push) are not only short but also alike to each other. Original BoWs model does not have the ability to capture their differences. Our approach can highlight their difference and solve this problem in an extent. On Rochester dataset, the results of “lookup in phone book”, “peel banana” and “use silverware” are unsatisfactory. Our approach can decrease those errors since temporal structure and salient parts are especially considered.

TABLE I compares the classification accuracy of our approach with state-of-the-arts on UT-Interaction and Rochester. Cluster number K is given to indicate the complexity of those approaches as there is no unified comparison approach. “K” is proportional to the algorithm’s dimensionality reduction ability. In our approach, Gradient descriptor [10] shows better results on UT-Interaction scene-2 and Rochester than 3D-SIFT because it extracts more

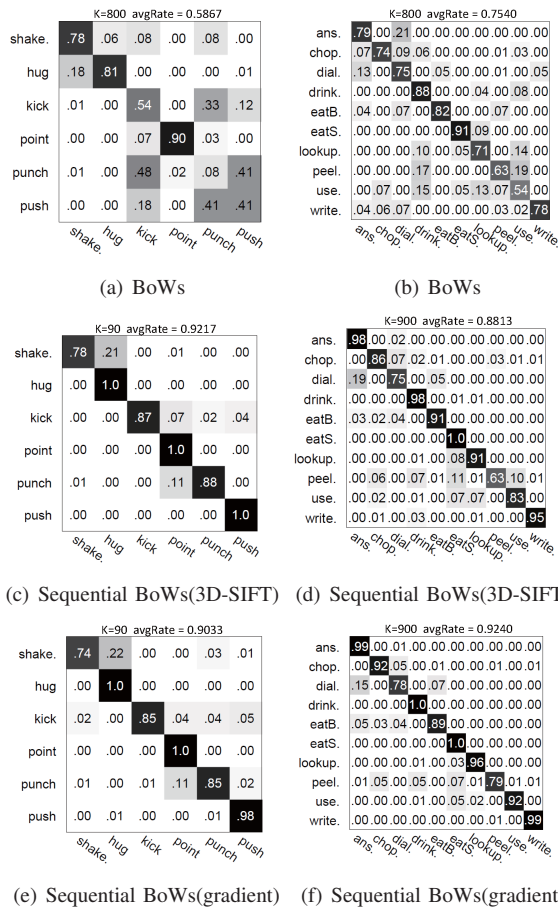


Fig. 7. Confusion matrices for scene-1 of UT-Interaction (left) and Rochester (right) is shown in column one and column two. From top to down, BoWs, sequential BoWs with 3D-SIFT descriptor and sequential BoWs with gradient descriptor. K is the cluster number.

feature points to describe the sub-actions more sufficiently. On UT-Interaction scene-1, our result is comparable to [6, 7], but manifest a faster computational speed. Noting that [6] focuses on the spatial structure between visual words, hence it can be combined with our approach. Approach in [7] aims at action prediction, and it conducts a iterate segment matching between classes and is inefficient to action classification. On both UT-Interaction scene-2 and Rochester datasets, our approach shows the best performance.

VI. CONCLUSIONS

In this paper, we present a novel approach called sequential BoWs for human action classification. It can reduce the ambiguities between classes sharing similar sub-actions. It captures the action's temporal sequential structure by segmenting the action into pieces. Salient pieces are stressed by assigning them higher weight and salience. Since our approach does not operate the descriptors directly, it can be combined with many mid-level descriptors. In experiments, the proposed approach is compared with the state-of-the-arts on two challenging datasets. Results show that our approach outperforms most existing BoWs based classification approaches especially on complex datasets with cluttered backgrounds and inter-class action ambiguities.

TABLE I
COMPARING PROPOSED APPROACH WITH STATE-OF-THE-ARTS

Method	scene1/K	scene2/K	Rochester/K
Dollar et al.[10]	58.67%/800	53.33%/800	75.40/800
Sun et al.[4]	82.67%/120	79.22%/120	-
Ryoo[7]	88%/800	77%/800	-
Satlin et al.[12]	-	-	80%/4000
Messing et al.[9]	-	-	89%/400
Liu et al.[6]	95%/450	86.67%/450	88%/500
Ours(3D-SIFT)	92.17%/90	85.83%/140	88.13%/900
Ours(gradient)	90.33%/90	91.83%/140	92.40%/900

REFERENCES

- [1] A. Kovashka, K. Grauman, "Learning a Hierarchy of Discriminative Space-Time Neighborhood Features for Human Action Recognition," IEEE Conference on Computer Vision and Pattern Recognition (CVPR), 2010, pp. 2046-2053.
- [2] J. Liu, M. Shah, "Learning Human Actions via Information Maximization," IEEE Conference on Computer Vision and Pattern Recognition (CVPR), 2008, pp. 1-8.
- [3] S. Savarese, A. Del Pozo, J.C. Niebles, L. Fei-Fei, "Spatial-Temporal Correlations for Unsupervised Action Classification," IEEE Workshop on Motion and Video Computing (WMVC), 2008, pp. 1-8.
- [4] Q. Sun and H. Liu, "Action Disambiguation Analysis Using Normalized Google-like Distance Correlogram," Asian Conference on Computer Vision (ACCV) 2012, 2013, pp. 425-437.
- [5] Q. Sun and H. Liu, "Learning Spatio-Temporal Co-occurrence Correlograms for Efficient Human Action Classification," IEEE International Conference on Image Processing (ICIP), 2013, pp. 3220-3224.
- [6] H. Liu, M. Liu, Q. Sun, "Learning Directional Co-occurrence for Human Action Classification," IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP), 2014.
- [7] M. S. Ryoo, "Human Activity Prediction: Early Recognition of Ongoing Activities from Streaming Videos," IEEE International Conference on Computer Vision (ICCV), 2011, pp. 1036-1043.
- [8] T. Glaser, L. Zelnik-Manor, "Incorporating Temporal Context in Bag-of-Words Models," IEEE International Conference on Computer Vision Workshops (ICCV Workshops), 2011, pp. 1562-1569.
- [9] R. Messing, C. Pal, and H. Kautz, "Activity Recognition Using the Velocity Histories of Tracked Keypoints," IEEE International Conference on Computer Vision (ICCV), 2009, pp. 104-111.
- [10] P. Dollar, V. Rabaud, G. Cottrell, and S. Belongie, "Behavior Recognition via Sparse Spatio-temporal Features," IEEE International Workshop on Visual Surveillance and Performance Evaluation of Tracking and Surveillance (VS-PETS), 2005, pp. 65-72.
- [11] G. Willems, T. Tuytelaars, L. V. Gool, "An Efficient Dense and Scale-Invariant Spatio-Temporal Interest Point Detector," European Conference on Computer Vision (ECCV), 2008, pp. 650-663.
- [12] M. S. Ryoo, C. C. Chen, J. K. Aggarwal, J. K. Aggarwal, and A. Roy-Chowdhury, "An Overview of Contest on Semantic Description of Human Activities," Recognizing Patterns in Signals, Speech, Images and Videos, 2010, pp. 270-285.
- [13] S. Satkin, M. Hebert, "Modeling the Temporal Extent of Actions," European Conference on Computer Vision (ECCV), 2010, pp. 536-548.
- [14] L. Shao, R. Mattivi, "Feature Detector and Descriptor Evaluation in Human Action Recognition," the ACM International Conference on Image and Video Retrieval, 2010, pp. 477-484.
- [15] R. Poppe, "A Survey on Vision-based Human Action Recognition," Image and Vision Computing, 2010, 28(6), pp. 976-990.
- [16] J. M. Chaqueta, E. J. Carmonaa, A. Fernandez-Caballero, "A Survey of Video Datasets for Human Action and Activity Recognition," Computer Vision and Image Understanding, 2013, 117(6), pp. 633-659.
- [17] P. Scovanner, S. Ali, M. Shah, "A 3-Dimensional SIFT Descriptor and its Application to Action Recognition," ACM Conference on Multimedia, 2007, pp. 357-360.