

Multi-Camera Action Dataset for Cross-Camera Action Recognition Benchmarking

Wenhui Li¹, Yongkang Wong², An-An Liu¹, Yang Li¹, Yu-Ting Su¹, Mohan Kankanhalli^{2,3}

¹School of Electronic Information Engineering, Tianjin University, China

²Interactive & Digital Media Institute, National University of Singapore, Singapore

³School of Computing, National University of Singapore, Singapore

Abstract

Action recognition has received increasing attention from the computer vision and machine learning communities in the last decade. To enable the study of this problem, there exist a vast number of action datasets, which are recorded under controlled laboratory settings, real-world surveillance environments, or crawled from the Internet. Apart from the “in-the-wild” datasets, the training and test split of conventional datasets often possess similar environments conditions, which leads to close to perfect performance on constrained datasets. In this paper, we introduce a new dataset, namely Multi-Camera Action Dataset (MCAD), which is designed to evaluate the open view classification problem under the surveillance environment. In total, MCAD contains 14,298 action samples from 18 action categories, which are performed by 20 subjects and independently recorded with 5 cameras. Inspired by the well received evaluation approach on the LFW dataset, we designed a standard evaluation protocol and benchmarked MCAD under several scenarios. The benchmark shows that while an average of 85% accuracy is achieved under the closed-view scenario, the performance suffers from a significant drop under the cross-view scenario. In the worst case scenario, the performance of 10-fold cross validation drops from 87.0% to 47.4%.

1. Introduction

Human action recognition has received increasing attention from the computer vision and machine learning community in the past few decades [8,17,21,27,31,40,41,46,51,

53,54]. Its importance is greatly driven by applications, such as human-computer interaction, action video indexing and retrieval, advanced video surveillance and so on.

In the early action recognition research, most of the research works were focused on the single-view learning problem. These works mainly focused on the extraction of robust feature representation (e.g. spatial features [17], spatio-temporal features [8,27], covariance descriptors [16,49], trajectories-based descriptor [46], etc.) and classification methodology [41]. More recently, semantic feature representations (i.e. local action attributes) were explored for improved action classification performance [21,31,48,54]. As the performance are saturating on the constrained datasets, several works have focused on cross-view learning problem [13,14,22] and cross-domain learning problem [5,7,9].

Cross-view learning aims to map features obtained from multiple views into a common feature space to handle the variations in visual appearance. In the case where a new action category is given, it can utilize the feature mapping model to perform action recognition between two different camera views. On the other hand, existing datasets often contain limited samples for each action category (see Table 1). To address this issue, cross-domain learning aims to leverage the small-scale data from target domain together with a large-scale data from an auxiliary domain to augment the generalization ability for model learning [45].

In the existing literature, many datasets are often collected under single camera view [15,34] or multiple views with overlapped observation [29,30,50]. Hence, it is hard to systematically evaluate the robustness of action recognition algorithms on similar yet different backgrounds and captured environments. On the other hand, the samples from large scale action recognition dataset collected from the Internet, such as UCF101 [43], consists of complex action captured from dynamic background environments. This type of datasets is ideal for deep learning based ap-

©2017 IEEE. Personal use of this material is permitted. Permission from IEEE must be obtained for all other uses, in any current or future media, including reprinting/republishing this material for advertising or promotional purposes, creating new collective works, for resale or redistribution to servers or lists, or reuse of any copyrighted component of this work in other works.

proaches [51,53].

Based on the above discussions, it is timely to have independently recorded multi-view constrained datasets, which provide standardized evaluation configuration to analyze the robustness of an action recognition system under unseen views. In this paper, we present a new Multi-Camera Action Dataset (MCAD), which consists of actions recorded with two types of CCTV cameras. Each camera has similar but slightly different FOVs, view perspective, image resolution, and background. The actions were independently performed on each camera view. Benchmark performance with single-view state-of-the-art algorithms indicate that this dataset is very challenging, especially for micro actions (*i.e.* action with small amount of motion area) and the cross-view action recognition scenario.

The rest of the paper is organized as follows. Section 2 reviews the existing datasets. Section 3 delineates the details of the proposed MCAD, where the benchmark is discussed in Section 4. Section 5 concludes the paper.

2. Dataset Review

2.1. Constrained Datasets

The constrained datasets are captured under controlled environments with constant background. Most of them were recorded under the indoor environment, which exhibited stable illumination conditions, fixed distance between person and cameras, and fixed direction of the actions.

The Weizmann dataset [15] contains clean and static background, and the participants perform actions around a small area. The KTH dataset [41] (see Fig. 1) is considered more challenging than the Weizmann dataset. It contains image sequences of human actions taken over from four scenarios and dynamic zoom variations. The dataset consists of relatively simple actions, such as “walking” and “jump”, with limited action variations. Literature has reported close to perfect performance on these datasets. Specifically, 100% classification accuracy on several action classes are reported in Weizmann dataset [15]. Different from these actions, there exist some datasets that recorded more complex actions. In the Activity of Daily Living (ADL) dataset [34], each activity is performed three times by five individuals of different shapes, size, gender, and ethnicity. Similarly, the TUM Breakfast dataset [25] comprises of actions related to breakfast preparation in various kitchens.

As the performance on these databases is saturating, several cross-view action recognition datasets were proposed. The first multi-view human action dataset is the INRIA Xmas Motion Acquisition Sequences (IXMAS) dataset [50] (see Fig. 1), which contains actions taken from 5 calibrated and synchronized cameras (4 side views and 1 top view). Subsequently, the Multicamera Human Action Video

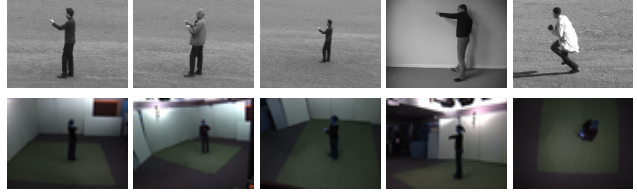


Figure 1: Sample images from constrained dataset. Top row: KTH dataset [41]; Bottom row: IXMAS dataset [50].

(MuHAVi) dataset [42] collected multiple primitive actions video data using 8 CCTV cameras located at 4 sides and 4 corners of a rectangular platform. Benefiting from the advances in depth sensing, the MV-TJU dataset [29] contains actions performed in both light and dark environment from two different cameras. Similarly, the Multi-modal & Multi-view & Interactive (M²I) dataset [30] extends the MV-TJU dataset by including person-person and person-object interactive action. Both the MV-TJU and M²I dataset consist of RGB image sequence, depth data and 3D skeleton data.

As many reported results on the constrained datasets are very good, these datasets are no longer regarded as challenging datasets for the action recognition problem. Furthermore, we argue that the actions are too simple when compared to the real world scenario. The action samples in these datasets are synchronized in all cameras, where the corresponding pairs have the same periodic properties. Several works are using this information to study the cross-view learning problem [55] and cross-domain learning problem [5,9,10,20]. In addition, we note that the camera views employed in the training stage are unlikely to have direct relationship (*i.e.* same view or overlapped region) with the test camera, especially for the surveillance application.

2.2. Consumer generated Datasets

The datasets of this category are generated by consumers and collected from the Internet, movies or personal video collections. These datasets are very challenging when compared with constrained datasets, due to its diversity in visual content, background complexity, and dynamic camera motion. Example of these datasets are shown in Fig. 2.

University of Central Florida (UCF) has collected several challenging human action datasets. UCF11 [32], UCF50 [38], and UCF101 [43] contain realistic videos and personal video collections collected from YouTube with different numbers of action classes. UCF Sports Action [38] consists of a set of actions in sports collected from a wide range of stock footage websites, including BBC Motion gallery and GettyImages. Other similar datasets include the Olympic sports dataset [35]. Moreover, the Human Motion Database (HMDB) [26] includes distinct action categories extracted from a wide range of sources. The Hollywood dataset [28] and the Hollywood2 dataset [33] con-

Table 1: Overview of existing action recognition dataset. ✕ indicates that the camera views are partially overlapped.

Dataset Type	Dataset Name	No. of Subjects	No. of Actions	No. of Views	No. of Samples	Image Resolution	Single-View	Multi-View	Indoor	Outdoor
Constrained	ADL [34]	5	10	1	150	$1,280 \times 720$	✓		✓	
	IXMAS [50]	10	11	5	1,650	320×240		✓	✓	
	KTH [41]	25	6	4	600	160×120	✓		✓	✓
	MuHAVi [42]	14	17	8	1,904	704×576		✓	✓	
	MV-TJU [29]	20	22	2	600	640×480		✓	✓	
	M ² I [30]	20	22	2	1,784	320×240		✓	✓	
	TUM Breakfast [25]	-	10	-	1,989	320×240	✓		✓	
	Weizmann [15]	9	10	1	90	180×144	✓			✓
Consumer	ASLAN [24]	-	8	-	233	-	✓		✓	✓
	HMDB [26]	-	51	-	6,766	-	✓		✓	✓
	Hollywood [28]	-	8	-	233	-	✓		✓	✓
	Hollywood2 [33]	-	12	-	3,669	-	✓		✓	✓
	Olympic sports [35]	-	16	-	800	-	✓		✓	✓
	Stanford 40 [52]	-	12	-	3,669	-	✓		✓	✓
	UCF11 [32]	-	11	-	3,040	-	✓		✓	✓
	UCF50 [38]	-	50	-	6,676	-	✓		✓	✓
	UCF101 [43]	-	101	-	13,320	-	✓		✓	✓
	UCF Sports [38]	-	10	-	184	720×480	✓		✓	✓
Surveillance	MSR [23]	10	3	2	-	320×240	✓		✓	✓
	iLIDS [36]	-	7	5	-	720×576		✕	✓	
	UCF Aerial [1]	-	9	-	-	-	✓			✓
	UCF-ARG [2]	12	10	3	1,440	$1,920 \times 1,080$		✓		✓
	UT-Interaction [39]	6	3	-	160	720×480	✓			✓
	MCAD (Proposed)	20	18	5	14,298	$1,280 \times 960$ & 704×576		✕	✓	

tain human actions distributed in the movies, which enable the comprehensive benchmark for human action recognition in the realistic and challenging settings. The Stanford 40 Action Dataset [52] contains images of humans performing 40 actions. Different from these datasets which are designed for action classification problem, Kliper-Gross *et al.* [24] proposed the Action Similarity Labeling (ASLAN) Challenge which contains 3697 action samples from 1571 unique YouTube videos divided into 432 non-trivial action categories. This benchmark focuses on the action verification problem.

2.3. Surveillance Datasets

The dataset of this category is captured with fixed view cameras under the real-world surveillance environments, which contains image sequences with complex background [36], aerial view [1], and crowded unconstrained environment [36].

The UCF Aerial Action dataset [1] was obtained using a R/C-controlled blimp equipped with an HD camera mounted on a gimbal. The collection represents a di-

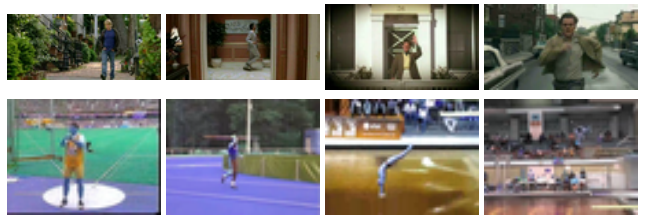


Figure 2: Sample images for consumer generated dataset. Top row: HMDB dataset [26]; Bottom row: Olympic sports dataset [35].

verse pool of actions featured at different heights and various viewpoints. The UT-Interaction dataset [39] focuses on human-human interactions in realistic environments in which each video contains at least one execution per interaction. The MSR dataset [23] was created in 2009 to study the behavior recognition algorithms in presence of clutter and dynamic backgrounds and other types of action variations. All the video sequences in this dataset are captured with clutter and moving backgrounds. The UCF-ARG dataset [2] is a multi-view real-world dataset which consists

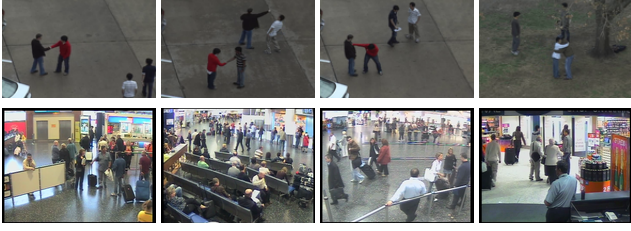


Figure 3: Sample images for surveillance dataset. Top row: UT-Interaction dataset [39]; Bottom row: iLIDS dataset [36].

of a ground camera, a rooftop camera, and an aerial camera mounted onto the payload platform of a helium balloon. The iLIDS dataset [36] is another multi-view real-world surveillance video in a busy airport. This dataset is also used in the TRECVID Surveillance Event Detection (SED) evaluation since 2008, where the presented action class remain challenging for the state-of-the-art approaches [36].

3. Multi-Camera Action Dataset

In this section, we delineate the details of the proposed dataset, namely Multi-Camera Action Dataset (MCAD)¹.

3.1. List of Recorded Actions

The MCAD consists of 9 single person daily actions and 9 person-object actions. These action categories are inherited from the KTH [41], IXMAS [50], and iLIDS [36] datasets. The action list and respective definition of each action are shown in Table 2. Among these actions, there are 7 actions that contains action with small amount of motion area², we denoted these actions as micro action. As demonstrated in Section 4, these micro actions are more challenging, especially the person-object actions.

In this dataset, we recruited a total of 20 human subjects. Each candidate repeats each action for 8 times (4 times during the day and 4 times in the evening) under one camera view. Different from multi-view datasets such as IXMAS [50] and MuHAVI [42] where several cameras are deployed to record an action sample synchronously, we use five cameras to record each action sample separately. Therefore, an algorithm designed for cross-view learning problem that deliberately explores the properties across two simultaneously recorded action is not applicable.

During the recording stage, we showed the subjects the list of actions and invited them to act freely with their personal preference. As a result, not only we observed high intra action class variation among different action samples, we also noticed some individuals acted differently across different camera view or section (*i.e.* daytime or nighttime).

¹ available via <http://mmas.comp.nus.edu.sg/MCAD/MCAD.html>

² action ID: {01, 02, 05, 10, 11, 12, & 13}

Table 2: List of actions and descriptions the proposed MCAD. Rows with RED and BLUE background color indicate single-person action and person-object action, respectively.

ActionID	Action Name	Action Description
01	Point	Someone points
02	Wave	Someone waves hand to catch peoples' attention
03	Jump	Someone jumps
04	Crouch	Someone crouches then stands up
05	Sneeze	Someone sneezes
06	SitDown	Someone sits down on a chair
07	StandUp	Someone stands up from a chair
08	Walk	Someone walks normally
09	PersonRun	Someone runs
10	CellToEar	Someone puts a cell phone to his/her ear
11	UseCellphone	Someone uses the cellphone to access information
12	DrinkingWater	Someone uses a bottle to drink water
13	TakePicture	Someone takes photos by using cellphone
14	ObjectGet	Someone bends or crouches to pick an object
15	ObjectPut	Someone puts down an object when walking
16	ObjectLeft	Someone walks and drops an object in this process
17	ObjectCarry	Someone walks with a bag
18	ObjectThrow	Someone throws a box to other place

For example, the *Jump* action in Fig. 4 demonstrates different posture on 5 randomly selected individuals. Under all recordings, the individuals were allowed to face any direction within cameras' FOV. This results in observable scale difference within the same camera view. The only exception is PTZ06 where the corresponding FOV is narrower than other camera views.

3.2. Environment Configuration

The MCAD is recorded with five unique cameras, including three static cameras (*i.e.* Cam04, Cam05 & Cam06) with fish eye effect and two Pan-Tilt-Zoom (PTZ) cameras (*i.e.* PTZ04 & PTZ06). These camera are mounted in a real-world surveillance environment. Among these cameras, the Cam04-PTZ04 and Cam06-PTZ06 pairs covered the same region with different FOV. The static camera has a resolution of 1280×960 pixels. The PTZ camera has a smaller FOV compared to the static camera, where the image resolution is 704×576 pixels.

The recording is carried out during both daytime and nighttime. In all cases, though the actions are independently recorded for each camera, the illumination condition is constant for each individual. However, due to the difference in the visual sensor and lens, we noticed observable difference in each camera view. For example, the *ObjectThrow* samples showed in Fig. 4 are all recorded during nighttime. Although the lighting conditions are the same for all cameras, the recorded footage on PTZ06 appears to be darker than Cam06, where both cameras observed the same region.

3.3. Evaluation Metric

In order to enable streamlined comparisons for future studies, we adopt the evaluation protocol from the Labeled Faces in the Wild (LFW) dataset [19]. The MCAD is di-



Figure 4: Sample images for the proposed MCAD dataset. Each row indicates unique action recorded with different individuals on 5 distinct camera views. Row 1 & 2 show samples recorded during day time where images in row 3 & 4 are recorded during night time.

vided into two sets, *i.e.* the Development Set and the Evaluation Set, The Development Set is recommended for parameters tuning. It consists of 10 randomly selected subjects from MCAD. In this work, we use the Leave-One-Subject-Out Cross Validation (LOSOVCV) strategy to evaluate the performance of an algorithm with various parameters. The optimal parameters are then applied to the evaluation set for reporting results. This protocol saves time during the comprehensive parameter search stage and creates an impartial condition for algorithm evaluation.

The Evaluation Set randomly divides all the subjects in MCAD into 10 training-test split³. For each training-test split, 12 subjects are selected as training data and the remaining 8 subjects as test set. We report the final 10-fold cross validation result with estimated mean accuracy and the standard error of the mean as in [19]. Specifically, the estimated mean accuracy $\hat{\mu}$ is given by

$$\hat{\mu} = \frac{\sum_{i=1}^{10} p_i}{10} \quad (1)$$

where p_i is the accuracy from i -th fold. The standard error

of the mean is given as

$$S_E = \frac{\hat{\sigma}}{\sqrt{10}} \quad (2)$$

where $\hat{\sigma}$ is the estimate of the standard deviation, given by

$$\hat{\sigma} = \sqrt{\frac{\sum_{i=1}^{10} (p_i - \hat{\mu})^2}{9}} \quad (3)$$

4. Benchmark

4.1. Baseline Algorithms

In this work, we benchmark the Bag-of-Words descriptor based method with four spatial-temporal local features and three encoding methods. Specifically, we selected Spatio-Temporal Interest Point (STIP) feature [27], Cuboid feature [8], Covariance matrices (denoted as Cov) [12], and Improved Dense Trajectory (IDT) [47]. For Cuboid feature, we use the parameter $\sigma = 2$ and $\tau = 1.5$ to extract up to 200 Cuboids from each action video, followed by Principle Component Analysis (PCA) [18] to reduce the dimensionality of the extracted feature to 100. For the covariance matrices, we first extract the 72-dimension HOF feature from Dense Trajectory (DT) feature to generate 72×72 dimensional covariance matrices X . Following [12] we compute

³ NOTE: The data split is available from the MCAD website

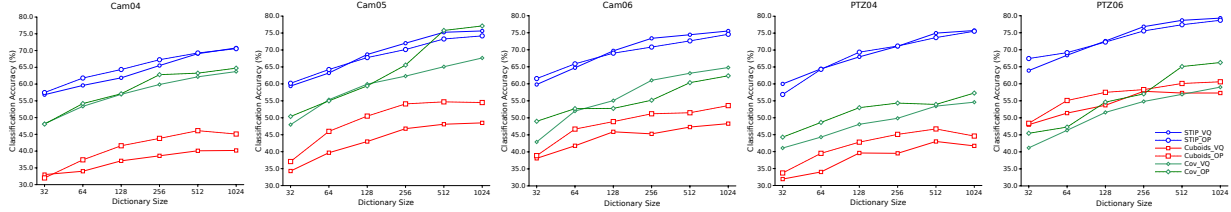


Figure 5: Performance of baseline algorithm on Development Set across various codebook size.

Table 3: Closed-view classification accuracy (%) on the Evaluation set. Cells with BLUE background color indicates the experiments are conducted with open-set classification constraints.

	STIP		Cuboids		Cov		IDT	
Training & Test Data	VQ	OP	VQ	OP	VQ	OP	FV	
Cam04	69.2 ± 0.6	69.2 ± 0.7	40.9 ± 0.6	40.5 ± 3.8	61.3 ± 0.4	63.5 ± 0.6	88.6 ± 0.5	83.5 ± 0.6
Cam05	74.0 ± 0.6	73.6 ± 0.7	47.7 ± 0.4	51.6 ± 0.5	65.4 ± 0.6	65.5 ± 0.8	91.6 ± 0.3	86.2 ± 0.4
Cam06	72.5 ± 0.4	72.9 ± 0.4	48.0 ± 0.7	51.7 ± 0.6	63.3 ± 0.4	60.2 ± 0.7	90.1 ± 0.3	83.6 ± 0.9
PTZ04	73.3 ± 0.4	73.4 ± 0.5	41.1 ± 0.4	45.5 ± 0.5	52.1 ± 0.5	55.2 ± 0.5	91.3 ± 0.3	86.5 ± 0.3
PTZ06	77.1 ± 0.6	76.3 ± 0.8	54.9 ± 0.8	57.8 ± 0.8	56.9 ± 0.5	57.3 ± 0.6	91.3 ± 0.7	87.0 ± 0.6
All Static Cameras	74.6 ± 0.4	74.9 ± 0.4	48.7 ± 0.3	52.4 ± 0.4	65.1 ± 0.4	68.5 ± 0.4	92.8 ± 0.2	84.2 ± 0.9
All PTZ Cameras	75.6 ± 0.5	76.1 ± 0.5	48.0 ± 0.6	52.3 ± 0.7	54.7 ± 0.5	58.4 ± 0.6	92.3 ± 0.2	87.2 ± 0.5
All Cameras	75.6 ± 0.4	76.2 ± 0.5	48.8 ± 0.4	53.7 ± 0.6	61.8 ± 0.4	66.1 ± 0.3	93.4 ± 0.4	84.2 ± 0.6

the Log-Euclidean vector representation of each X and use this representation as covariance feature.

In our first set of baseline methods, we adopt two encoding methods for STIP, Cuboid and Cov features. In the first encoding method, we utilized Kmeans++ [4] to learn the codebook and use Vector Quantization to encode each local feature, followed by mean pooling to generate the descriptor. We denote this baseline method as *FeatureName_VQ*. For the second encoding method, we adopt a sparse coding approach, where the K-SVD algorithm [3] is utilized for codebook learning and Orthogonal Matching Pursuit (OMP) algorithm [44] is used for encoding. This baseline is denoted as *FeatureName_OP*. For the IDT feature, we follow [47] and use Fisher Vector (FV) encoding to generate the descriptor. Specifically, we first use Gaussian Mixture Model (GMM) to learn the codebook. Unlike VQ and OMP encoding method, FV encodes both the first and second order statistics between the video descriptors and a GMM. We denote this baseline as IDT_FV.

Given a descriptor, we used an SVM with χ^2 kernel for classification. The dimensionality of IDT_FV descriptor is too high for χ^2 kernel SVM. We select a linear SVM, which shows good results in classification for descriptor with high dimensionality [11,37]. The aforementioned classification is deployed for closed-set classification scenario. In this benchmark, we also employ an open-set linear SVM classifier [6] to evaluate the performance under open-set scenario. Based on the preliminary experiment, near and far plane pressures is fixed to 0.4 and 1.0, respectively.

4.2. Evaluation under Closed View Scenario

In this section, we evaluate the benchmark performance with closed view recognition scenario, *i.e.* the camera view of the test data is the same as that for the training data. Specifically, we restrict the training and test data from same camera source, while the subjects can only appear in either the training or test data. Three types of camera source scenarios are evaluated, namely Single Camera (SC), Same camera Type (ST), and All Cameras (AC).

First, we use the Development Set to fine-tune the optimal codebook size and parameters of SVM classifier under closed-set classification scenario. Fig. 5 shows the performance on all 5 camera views. Across all camera views, the accuracy gradually increases and saturates when the codebook size is set to 1024. The only exception is for the Cuboids feature based descriptor, which the performance under codebook size of 512 is the best. Due to the limits of computational resources, the codebook size of IDT_FV is evaluated up to 256. Based on the optimal performance, we conducted 10-fold cross validation on the Evaluation Set. The results are shown in Table 3 and the category wise performance of the best performing descriptor are shown in Fig. 6. The key findings are as follows:

1. IDT_FV consistently achieved the highest mean accuracy and lowest standard error on all scenarios, which is consistent with the reported performance on other datasets [47]. For the STIP and Cuboids features, we notice that the performance with PTZ06 is better than

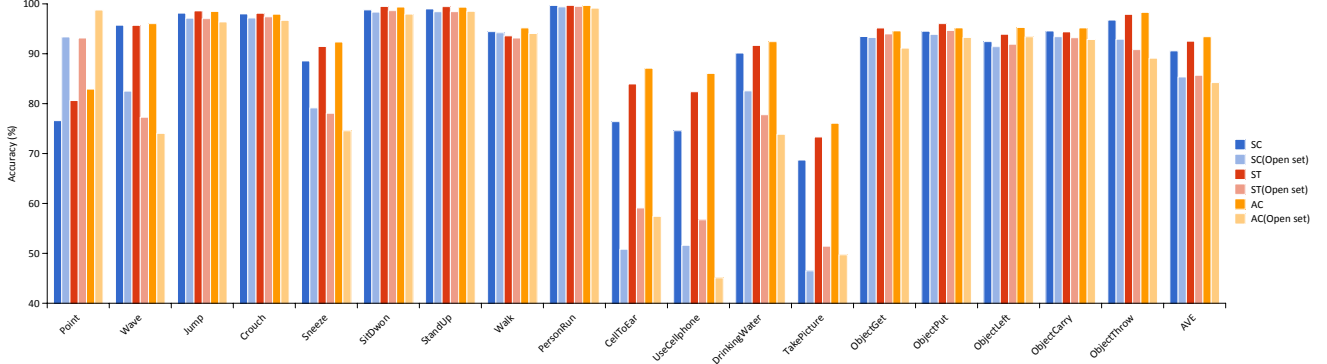


Figure 6: Mean classification accuracy of IDT_FV descriptor on the Evaluation set under Single Camera (SC), Same Type (ST), and All Camera (AC) scenarios.

for other camera views. As shown in Fig. 4, the FOV of PTZ06 is narrower and the size of each person is more consistent than in other cameras. Hence, the extracted local features is more consistent.

2. We observed that when the available training data increases from single camera view to all cameras, the performance of Cov_OP and IDT_FV increases. However, this is not true for the VQ based encoding method, where some of the single camera view scenarios report the best performance. This suggests that IDT_FV not only is more discriminative than other baseline methods, it also is more robust in handling training data with large environmental variation.
3. From Fig. 6, we found that the performance of micro actions, such as *Sneeze*, *CellToEar*, and *TakePicture*, is significant lower than the other class of actions. One reasons is that the available local features on the micro actions are fewer. On the other hand, action with large spatial movement, such as *PersonRun* and *Jump*, can be accurately recognized. In generally, person-object actions are slightly harder to recognize than single person actions.

In this work, we also evaluate the performance of the baseline methods under open-set scenario. Specifically, the classifier is required to identify whether the given sample belongs to one of the known class given in the training set or rejects it as an unknown sample, which is a realistic scenario in real-world applications. Based on the findings on the above section, we only show the performance of IDT_FV approach. From Table 3 (in BLUE background color) and Fig. 6, the performance of IDT_FV degraded on all scenarios. In our preliminary experiment, we also evaluated the performance of open-set action recognition with KTH [41], M²I [30], and IXMAS [50], where similar performance trend are observed. Under MCAD, the most significant performance drop is with the *UseCellPhone* ac-

tion, where the performance under all camera cases dropped from 86.7% to 45.1%.

Surprisingly, the performance of *Point* action increased for SC, ST, and AC scenarios. To investigate this, we examine the confusion matrix of both closed-set and open-set scenario under AC scenario. As shown in Fig. 7, we find that the *Point* action is easier to confuse with other actions in the closed-set classification scenario, where the probability of *Point* action to be misclassified as *CellToEar* action and *TakePicture* action are 0.07 and 0.04, respectively. While in the open-set scenario, the probability of *TakePicture* action reduced to 0.01. We also observed that the performance of the micro actions with object (*i.e.* *CellToEar*, *UseCellphone*, *DrinkingWater*, and *TakePicture*) suffers significantly. If we closely examine the confusion matrix, most of the test samples under these actions are mostly misclassified as *Point* action. Similar to the *Point* action, these action classes contain the arm motion action, which might make it harder to distinguish with interest point based descriptor under the open-set scenario.

4.3. Evaluation with Open View Classification

In this section, we evaluate the benchmark performance for the open view recognition scenario, *i.e.* the camera view of the test data has never been seen in the training phase. We comprehensively evaluate all single camera cross view classification cases, where data from one camera is selected to initialize the codebook and train the classifier (*i.e.* source view), and the evaluation is conducted on the selected camera (*i.e.* target view). For all cases, the subjects in each data split is identical to those in Section 4.2. We applied the optimal codebook size from the previous section and the SVM parameters are fine-tuned on the Development Set. Both the close-set and open-set classification scenarios are evaluated. The results with IDT_FV on MCAD are shown in Table 4. Furthermore, we also conducted the same experiment on the synchronous IXMAS dataset (see Table 5).

Table 4: Open-view classification accuracy (%) with IDT_FV on the Evaluation set. First column indicates the source of the training data while the remaining columns are the evaluation with respective test image sequences. The diagonal entries (*i.e.* cells with RED background color) are the classification accuracy of closed view SC scenario for comparison purposes.

	Closed-Set Classification					Open-Set Classification				
Training Data	Cam04	Cam05	Cam06	PTZ04	PTZ06	Cam04	Cam05	Cam06	PTZ04	PTZ06
Cam04	88.6 ± 0.6	81.5 ± 0.6	75.6 ± 0.6	74.6 ± 0.7	63.4 ± 0.7	83.5 ± 0.6	73.1 ± 1.0	64.1 ± 1.1	65.8 ± 0.9	52.9 ± 1.3
Cam05	80.9 ± 0.7	91.6 ± 0.3	71.6 ± 0.7	73.9 ± 0.4	62.3 ± 0.8	75.2 ± 0.8	86.2 ± 0.4	59.8 ± 0.8	64.3 ± 0.5	51.6 ± 1.2
Cam06	73.0 ± 0.6	72.8 ± 0.4	90.1 ± 0.3	65.8 ± 0.4	68.5 ± 0.5	68.0 ± 0.7	65.1 ± 0.9	83.6 ± 0.9	58.7 ± 0.7	58.1 ± 1.2
PTZ04	75.7 ± 0.4	72.9 ± 0.5	67.8 ± 0.9	91.3 ± 0.3	59.3 ± 1.1	68.5 ± 0.7	63.7 ± 0.7	54.0 ± 0.7	86.5 ± 0.3	47.4 ± 1.0
PTZ06	57.4 ± 0.6	59.1 ± 0.5	60.3 ± 0.4	56.9 ± 0.4	91.3 ± 0.7	53.0 ± 0.6	53.5 ± 0.5	54.0 ± 0.6	50.2 ± 0.6	87.0 ± 0.6

Table 5: Open-view classification accuracy (%) with IDT_FV on the Evaluation set in IXMAS dataset. First column indicates the source of the training data while the remaining columns are the evaluation with respective test image sequences. The diagonal entries (*i.e.* cells with RED background color) are the classification accuracy of closed view SC scenario for comparison purposes.

	Closed-Set Classification					Open-Set Classification				
Training Data	Cam0	Cam1	Cam2	Cam3	Cam4	Cam0	Cam1	Cam2	Cam3	Cam4
Cam0	95.6 ± 0.9	87.0 ± 2.2	48.1 ± 2.2	64.7 ± 2.6	-	88.6 ± 2.1	74.0 ± 2.2	47.0 ± 1.9	60.0 ± 2.9	-
Cam1	71.2 ± 2.2	95.8 ± 1.0	39.0 ± 2.2	43.6 ± 2.2	-	60.5 ± 4.2	94.8 ± 1.2	36.1 ± 3.3	36.6 ± 4.7	-
Cam2	71.2 ± 2.5	69.9 ± 1.8	94.3 ± 0.7	73.5 ± 1.7	-	58.4 ± 1.5	48.8 ± 4.0	90.7 ± 1.2	55.3 ± 2.3	-
Cam3	66.8 ± 3.4	72.5 ± 2.6	68.1 ± 1.7	93.5 ± 1.8	-	59.5 ± 3.3	51.4 ± 3.0	50.4 ± 1.8	90.4 ± 2.0	-
Cam4	-	-	-	-	94.3 ± 1.1	-	-	-	-	87.8 ± 1.7

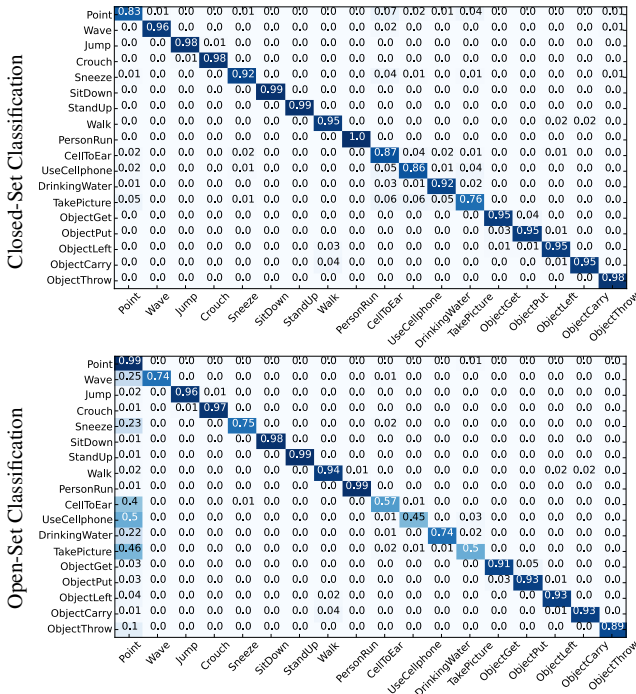


Figure 7: Confusion matrix of IDT.FV descriptor on the Evaluation set under AC scenario.

Overall, the performance of cross view recognition drops significantly when compared to the single camera case. The results are expected as the training and test data has significant difference in view perspective, FOV, image quality, and pixel resolution. Consistent with the previous section, the performance with open-set classification method further reduce the performance. We also observed that the performance is more stable when the evaluated camera view has

similar properties as that of the training data. For instance, the Cam04-Cam05 pair report an average of 74.15% on open-set classification scenario, where the corresponding performance on Cam04-PTZ04 pair is around 75.4%. The view conditions of Cam06 is significantly different from the other cameras, and registers worst performance on all cross-view evaluation. In Table 5, we deliberately do not perform cross view evaluation for Cam4 because Cam4 is a top view camera (see Fig. 1) and does not exhibit visually favorable properties for meaningful action recognition task.

Finally, we highlight that the open view evaluation is essential to assess the robustness of any proposed algorithm. Different from action recognition with consumer generated data (*e.g.* egocentric video or crowdsourced dataset), it is impractical to collect surveillance video data from all possible conditions for training purposes. It is important to point out that for dataset that are synchronously recorded (*e.g.* IXMAS dataset), the evaluation needs to carefully designed such that temporal self-similarity is not utilized to improve the performance. In our future work, we plan to evaluate view-invariant action recognition algorithms on the open view scenario.

5. Conclusion

In this paper, we presented a new action recognition dataset, namely Multi-Camera Action Dataset (MCAD), which is designed to evaluate the open view action classification problem. Different from existing multi-view datasets, the samples in MCAD are independently recorded with 5 cameras and 20 subjects, and contains a total of 14,298 action samples. Inspired by the LFW dataset, we designed a standard evaluation protocol and benchmarked

MCAD under several scenarios.

Acknowledgment

This research is supported by the National Research Foundation, Prime Ministers Office, Singapore under its International Research Centre in Singapore Funding Initiative.

References

- [1] UCF aerial action dataset. <http://server.cs.ucf.edu/vision/aerial/index.html>. 3
- [2] UCF aerial camera, rooftop camera and ground camera dataset. <http://crcv.ucf.edu/data/UCF-ARG.php>. 3
- [3] M. Aharon, M. Elad, and A. Bruckstein. K-SVD: An algorithm for designing overcomplete dictionaries for sparse representation. *IEEE Transactions on signal processing*, 54(11):4311–4322, 2006. 6
- [4] D. Arthur and S. Vassilvitskii. k-means++: The advantages of careful seeding. In *ACM-SIAM symposium on Discrete algorithms*, pages 1027–1035, 2007. 6
- [5] F. R. Bach, G. R. G. Lanckriet, and M. I. Jordan. Multiple kernel learning, conic duality, and the SMO algorithm. In *ICML*, 2004. 1, 2
- [6] A. Bendale and T. Boulton. Towards open world recognition. In *CVPR*, pages 1893–1902, 2015. 6
- [7] Z. Cui, W. Li, D. Xu, S. Shan, X. Chen, and X. Li. Flowing on Riemannian manifold: Domain adaptation by shifting covariance. *IEEE Transactions on Cybernetics*, 44(12):2264–2273, 2014. 1
- [8] P. Dollár, V. Rabaud, G. Cottrell, and S. Belongie. Behavior recognition via sparse spatio-temporal features. In *IEEE International Workshop on VS-PETS*, pages 65–72, 2005. 1, 5
- [9] L. Duan, I. W. Tsang, and D. Xu. Domain transfer multiple kernel learning. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 34(3):465–479, 2012. 1, 2
- [10] L. Duan, D. Xu, I. W. Tsang, and J. Luo. Visual event recognition in videos by learning from web data. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 34(9):1667–1680, 2012. 2
- [11] R.-E. Fan, K.-W. Chang, C.-J. Hsieh, X.-R. Wang, and C.-J. Lin. LIBLINEAR: A library for large linear classification. *Journal of Machine Learning Research*, 9:1871–1874, 2008. 6
- [12] M. Faraki, M. Palhang, and C. Sanderson. Log-Euclidean bag of words for human action recognition. *IET Computer Vision*, 9(3):331–339, 2014. 5
- [13] A. Farhadi and M. K. Tabrizi. Learning to recognize activities from the wrong view point. *Lecture Notes in Computer Science*, 5302:154–166, 2008. 1
- [14] Z. Gao, W. Nie, A. Liu, and H. Zhang. Evaluation of local spatial-temporal features for cross-view action recognition. *Neurocomputing*, 173:110–117, 2016. 1
- [15] L. Gorelick, M. Blank, E. Shechtman, M. Irani, and R. Basri. Actions as space-time shapes. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 29(12):2247–2253, 2007. 1, 2, 3
- [16] L. Han, W. Liang, X. Wu, and Y. Jia. Human action recognition using discriminative models in the learned hierarchical manifold space. In *AFGR*, pages 1–6, 2008. 1
- [17] C. Harris and M. Stephens. A combined corner and edge detector. In *Alvey Vision Conference*, pages 23.1–23.6, 1988. 1
- [18] T. Hastie, R. Tibshirani, and J. Friedman. *The Elements of Statistical Learning: Data Mining, Inference, and Prediction*, volume 1. Springer, 2001. 5
- [19] G. B. Huang, M. Ramesh, T. Berg, and E. Learned-Miller. Labeled faces in the wild: A database for studying face recognition in unconstrained environments. Technical report, Technical Report 07-49, University of Massachusetts, Amherst, 2007. 4, 5
- [20] H. D. III. Frustratingly easy domain adaptation. In *ACL*, 2007. 2
- [21] A. Jain, A. Gupta, M. Rodriguez, and L. S. Davis. Representing videos using mid-level discriminative patches. In *CVPR*, pages 2571–2578, 2013. 1
- [22] Y. Jiang, F. Chung, S. Wang, Z. Deng, J. Wang, and P. Qian. Collaborative fuzzy clustering from multiple weighted views. *IEEE Transactions on Cybernetics*, 45(4):688–701, 2015. 1
- [23] Y. Ke, R. Sukthankar, and M. Hebert. Event detection in crowded videos. In *ICCV*, pages 1–8, 2007. 3
- [24] O. Kliper-Gross, T. Hassner, and L. Wolf. The action similarity labeling challenge. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 34(3):615–621, 2012. 3
- [25] H. Kuehne, A. Arslan, and T. Serre. The language of actions: Recovering the syntax and semantics of goal-directed human activities. In *CVPR*, pages 780–787, 2014. 2, 3
- [26] H. Kuehne, H. Jhuang, E. Garrote, T. Poggio, and T. Serre. HMDB: a large video database for human motion recognition. In *ICCV*, pages 2556–2563, 2011. 2, 3
- [27] I. Laptev and T. Lindeberg. Space-time interest points. In *ICCV*, pages 432–439, 2003. 1, 5
- [28] I. Laptev, M. Marszalek, C. Schmid, and B. Rozenfeld. Learning realistic human actions from movies. In *CVPR*, pages 1–8, 2008. 2, 3
- [29] A.-A. Liu, Y.-T. Su, P.-P. Jia, Z. Gao, T. Hao, and Z.-X. Yang. Multiple/single-view human action recognition via part-induced multitask structural learning. *IEEE Transactions on Cybernetics*, 45(6):1194–1208, 2015. 1, 2, 3
- [30] A.-A. Liu, N. Xu, W.-Z. Nie, Y.-T. Su, Y. Wong, and M. Kankanhalli. Benchmarking a multi-modal & multi-view & interactive dataset for human action recognition. *IEEE Transactions on Cybernetics*, 2016. in press. 1, 2, 3, 7
- [31] J. Liu, B. Kuipers, and S. Savarese. Recognizing human actions by attributes. In *CVPR*, pages 3337–3344, 2011. 1

- [32] J. Liu, J. Luo, and M. Shah. Recognizing realistic actions from videos “in the wild”. In *CVPR*, pages 1996–2003, 2009. 2, 3
- [33] M. Marszalek, I. Laptev, and C. Schmid. Actions in context. In *CVPR*, pages 2929–2936, 2009. 2, 3
- [34] R. Messing, C. Pal, and H. Kautz. Activity recognition using the velocity histories of tracked keypoints. In *ICCV*, pages 104–111, 2009. 1, 2, 3
- [35] J. C. Niebles, C. Chen, and F. Li. Modeling temporal structure of decomposable motion segments for activity classification. *Lecture Notes in Computer Science*, 6312:392–405, 2010. 2, 3
- [36] P. Over, J. Fiscus, G. Sanders, D. Joy, M. Michel, G. Awad, A. Smeaton, W. Kraaij, and G. Quénot. TRECVID 2014 – an overview of the goals, tasks, data, evaluation mechanisms and metrics. In *Proceedings of TRECVID*, page 52, 2014. 3, 4
- [37] D. A. Reynolds. Gaussian mixture models. In *Encyclopedia of Biometrics*, pages 659–663. 2009. 6
- [38] M. D. Rodriguez, J. Ahmed, and M. Shah. Action MACH a spatio-temporal maximum average correlation height filter for action recognition. In *CVPR*, 2008. 2, 3
- [39] M. S. Ryoo and J. K. Aggarwal. Spatio-temporal relationship match: Video structure comparison for recognition of complex human activities. In *ICCV*, pages 1593–1600, 2009. 3, 4
- [40] M. Sapienza, F. Cuzzolin, and P. H. Torr. Learning discriminative space-time action parts from weakly labelled videos. *International Journal of Computer Vision*, 110(1):30–47, 2014. 1
- [41] C. Schuldt, I. Laptev, and B. Caputo. Recognizing human actions: a local SVM approach. In *ICPR*, pages 32–36, 2004. 1, 2, 3, 4, 7
- [42] S. Singh, S. A. Velastin, and H. Ragheb. MuHAVI: A multicamera human action video dataset for the evaluation of action recognition methods. In *AVSS*, pages 48–55, 2010. 2, 3, 4
- [43] K. Soomro, A. R. Zamir, and M. Shah. UCF101: A dataset of 101 human actions classes from videos in the wild. Technical report, November 2012. 1, 2, 3
- [44] J. A. Tropp and A. C. Gilbert. Signal recovery from random measurements via orthogonal matching pursuit. *IEEE Transactions on Information Theory*, 53(12):4655–4666, 2007. 6
- [45] G. Wang, F. Wang, T. Chen, D. Yeung, and F. H. Lochovsky. Solution path for manifold regularized semisupervised classification. *IEEE Transactions on Systems, Man, and Cybernetics, Part B*, 42(2):308–319, 2012. 1
- [46] H. Wang, A. Kläser, C. Schmid, and C.-L. Liu. Dense trajectories and motion boundary descriptors for action recognition. *International Journal of Computer Vision*, 103(1):60–79, 2013. 1
- [47] H. Wang and C. Schmid. Action recognition with improved trajectories. In *ICCV*, pages 3551–3558, 2013. 5, 6
- [48] L. Wang, Y. Qiao, and X. Tang. Mining motion atoms and phrases for complex action recognition. In *ICCV*, pages 2680–2687, 2013. 1
- [49] L. Wang and D. Suter. Learning and matching of dynamic shape manifolds for human action recognition. *IEEE Transactions on Image Processing*, 16(6):1646–1661, 2007. 1
- [50] D. Weinland, E. Boyer, and R. Ronfard. Action recognition from arbitrary views using 3D exemplars. In *ICCV*, pages 1–7, 2007. 1, 2, 3, 4, 7
- [51] Z. Xu, Y. Yang, and A. G. Hauptmann. A discriminative CNN video representation for event detection. In *CVPR*, pages 1798–1807, 2015. 1, 2
- [52] B. Yao, X. Jiang, A. Khosla, A. L. Lin, L. J. Guibas, and F. Li. Human action recognition by learning bases of action attributes and parts. In *ICCV*, pages 1331–1338, 2011. 3
- [53] S. Zha, F. Luisier, W. Andrews, N. Srivastava, and R. Salakhutdinov. Exploiting image-trained CNN architectures for unconstrained video classification. In *BMVC*, pages 60.1–60.13, 2015. 1, 2
- [54] J. Zhang, H. Lin, W. Nie, L. Chaisorn, Y. Wong, and M. S. Kankanhalli. Human action recognition bases on local action attributes. *JEET*, 10(3), 2015. 1
- [55] J. Zheng, Z. Jiang, P. J. Phillips, and R. Chellappa. Cross-view action recognition via a transferable dictionary pair. In *BMVC*, pages 1–11, 2012. 2