

Action Recognition in the Frequency Domain*

Anh Tran[†], Jinyan Guan[†], Thanima Pilantanakitti[‡], Paul Cohen[‡]

University of Arizona

[†]Department of Computer Science

[‡]School of Information: Science, Technology, and Arts

{trananh, jguan1, tpilanta, prcohen}@email.arizona.edu

Abstract

In this paper, we describe a simple strategy for mitigating variability in temporal data series by shifting focus onto long-term, frequency domain features that are less susceptible to variability. We apply this method to the human action recognition task and demonstrate how working in the frequency domain can yield good recognition features for commonly used optical flow and articulated pose features, which are highly sensitive to small differences in motion, viewpoint, dynamic backgrounds, occlusion and other sources of variability. We show how these frequency-based features can be used in combination with a simple forest classifier to achieve good and robust results on the popular KTH Actions dataset.

1 Introduction

Human action recognition research has attracted considerable attention in recent years due to its practical applications in areas such as video surveillance, robotics, human-computer interaction, video indexing and retrieval, scene understanding and analysis, behavioral biometrics, biomechanics, and others.

Typical recognition scenarios often include variations in motion, illumination and viewpoint, partial occlusions, variable execution rates and anthropometry of actors involved, changes in backgrounds, and so forth (Aggarwal and Cai 1999; Moeslund and Granum 2001). These conditions pose great challenges for researchers and often induce much variability in the data.

In this paper, we present a strategy to reduce the effects on classifier performance of some of these kinds of variability. Rather than working with data in the temporal domain, it is sometimes better to extract features of the data in the frequency domain. We apply this idea to two commonly used features in human action recognition research,

optical flow and articulated pose, and show how their corresponding frequency-domain features can be effectively used for classification on the KTH Actions dataset from Schuldt, Laptev, and Caputo (2004). We adopt the efficient randomized forest-based approximate nearest-neighbor approach presented by O’Hara and Draper (2012) and apply it to our frequency-domain features to build *frequency forest* classifiers.

Details of the frequency domain representation and classification method are described in Sections 3 and 4, respectively. Section 2 presents related work, and Section 5 describes our experimental procedures and results. Section 6 concludes this paper.

2 Related Work

There is a large body of literature on human action recognition research, detailing many innovative learning algorithms and novel representations for actions. Several informative surveys are available, including (Aggarwal and Cai 1999; Moeslund and Granum 2001; Gavrilu 1999; Wang and Singh 2003; Turaga et al. 2008). Here, we will briefly review the more commonly used and recently developed representations and classifiers of actions.

A popular approach is to use localized space-time features extracted from video sequences. Schuldt, Laptev, and Caputo (2004) demonstrated that local measurements in terms of spatio-temporal interest points (STIP) can sufficiently represent the complex motion patterns of various actions. Recent research also achieved success by representing *tracklets* of actions as points on Grassmann manifolds, where each track is modeled as a 3-dimensional data cube with axes of width, height, and frame number (O’Hara and Draper 2012; Lui, Beveridge, and Kirby 2010).

Asides from low-level representations, there are work that focus on higher level representations of actions. For instance, Kerr, Tran, and Cohen (2011) described each action sequence as a collection of propositions that have truth values over some time intervals. They showed that each action has a corresponding set of signature *fluents* — intervals during which propositions are true — that can be learned, and that finite state machine models of actions can be constructed from these signatures. Recently, Sadanand and Corso (2012) presented a new high-level representation of videos called Action Bank. Inspired by the *object bank*

*This work was partially supported by the DARPA Minds Eye program. The U.S. Government is authorized to reproduce and distribute reprints for Governmental purposes notwithstanding any copyright annotation thereon. The views and conclusions contained herein are those of the authors and should not be interpreted as necessarily representing the official policies or endorsements, either expressed or implied, of DARPA or the U.S. Government.

approach to image representation, an *action bank* representation comprises the collected output of many pre-trained template-based action detectors that each produce a correlation volume.

Actions can also be modeled in terms of the optical flow of various regions of the image. Yacoob and Black (1998) represented actions using the optical flow of different body parts (e.g., torso, thighs, calfs, feet, and arms). Danafar and Gheissari (2007) segmented the body into more coarse regions (e.g., head-neck, body-hands, and legs) and represented actions using histograms of flow in these regions.

Innovations in pose estimation technology have also inspired representations centered on features extracted from articulated body parts. In this area, some researchers like to work with articulated pose in 2D, as in (Sheikh, Sheikh, and Shah 2005; Lv and Nevatia 2007); while others prefer to avoid the challenges of 2-dimensional image data by directly recording joints coordinates in 3D using the increasingly accessible RGBD cameras or other commercial motion capture systems, as in (Campbell and Bobick 1995; Li, Zhang, and Liu 2010; Wang 2012; Sung et al. 2011; Xia, Chen, and Aggarwal 2012).

3 Representation: Frequency Features

Frequency features for action recognition are representations in the frequency domain that model the body motions associated with each action. For example, the frequency features for a *walk* action might capture the rhythmic swinging of the arms or the periodic movements of the legs, and not just the fact that a blob of energy shifted across the screen.

In this work, we assume that the body has been localized in the frame using some combinations of computer vision detection and tracking algorithms, or any other similar method. This localization problem is itself an open challenge in computer vision. However, it is not the focus of this paper, so the work reported here is done with hand-annotated videos in which bounding boxes have been drawn around the body. (Some evidence suggests that feature-based methods require localization to perform well on datasets that have dynamic backgrounds (Ryoo and Chen 2010), and that bounding-box style localization is easier to use than alternatives like silhouette or pose extraction.) Our bounding boxes were created using the VATIC annotation tool (Vondrick, Ramanan, and Patterson 2010).¹ Examples of the boxes can be seen in Figure 1.

The rest of this section explains how we represent actions in the frequency domain using optical flow and articulated pose features that are localized by bounding boxes.

Optical Flow We extract optical flow from images using the algorithm described in Liu (2009). Once extracted, we use the bounding box to localize the flow of the actor, and then further divide the box area into smaller regions. Following a similar process to Danafar and Gheissari (2007), we split each bounding box into 5 different subregions (see Figure 2): head, left torso/arm, right torso/arm, left leg, and

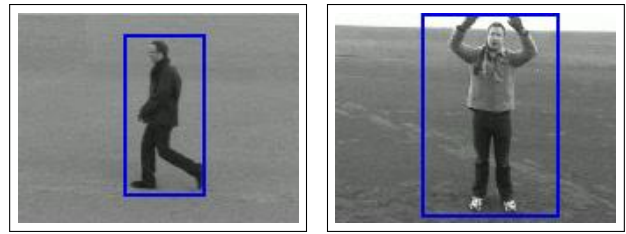


Figure 1: Example annotated bounding box localization for two different video frames.

right leg. The corresponding left and right regions are split evenly in the horizontal direction. The head region occupies $1/5$ of the bounding box's height, while the torso/arms and legs regions are each $2/5$ of the height.

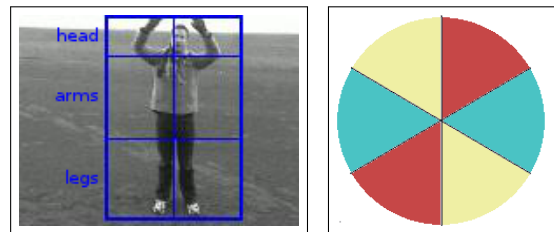


Figure 2: Illustration of optical flow partitioning for a bounding box. The 5 subregions representing different areas of the body are shown in the left image. The right image shows the 6 directional bins used to partition the flow within each body subregion.

The flow within each subregion is binned into six different directions, as in Figure 2. Within each bin, the proportion of vectors that fall into the bin with respect to all vectors in the subregion, as well as the average magnitude of all flow vectors in the bin are computed.

These calculations yield *optical flow features* such as the proportion of flow vectors falling into the rightward bin in the head subregion, or the average magnitude of flow vectors falling into the up-leftward bin in the right leg subregion. Optical flow features are calculated for each frame, but we are interested in their time series across all frames, and particularly in transformations of these time series into the frequency domain by the Fourier transform.

After the Fourier transformation, we examine the power spectrum of each series and take the first N components of the spectrum to be a frequency-domain feature. (In our experiments, $N = 25$.) Figure 3 shows some example time-series features and their corresponding power spectra. For short video segments that yield fewer than N components, we simply recycle the segment. This method allows us to convert the time series of each optical flow feature into a corresponding frequency feature, and, in fact, we can calculate frequency features for any combination of optical flow features. For example, among the 31 frequency features that we used for action recognition, we derived the right plus left flow proportions of each subregion, the upper-right plus lower-left proportions of each subregion, the lower-right

¹All annotations and relevant code are publicly available at: <https://code.google.com/p/ua-gesture/>

plus upper-left proportions of each subregion, the average flow magnitudes in different directions of each subregion, and the average flow of the entire bounding box.

Articulated Pose To extract articulated pose from images, we use the algorithm by Yang and Ramanan (2011) which returns all detected poses for each frame. We then find the pose that best matches each bounding box (i.e., the pose with the highest detection score that fits in the area of the bounding box) to generate tracks of articulated poses. We note that due to the low resolution of the data, some images needed to be enlarged to $2\times$ or $3\times$ magnification for the pose estimation algorithm to work properly.

We converted the 26-joint pose format given by Yang and Ramanan (2011) to a simpler 15-joint format that is more compatible with other pose datasets (see Figure 4). The time series for each joint is smoothed and all poses are then standardized to be within a unit square. This ensures that poses are all roughly equal in size, allowing some robustness to variability in the size of the figure in a frame.

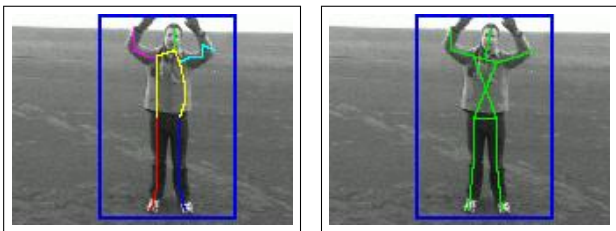


Figure 4: An example articulated pose output by Yang and Ramanan (2011) that best matches the bounding-box (left image) and the corresponding reformatted pose with fewer joints (right image).

As with the optical flow features, we treat the different measurements, or relationships, between joint trajectories as temporal data series and convert them into corresponding frequency spectra using the Fourier transform. We take the top $N = 25$ components of each associated power spectrum as a frequency feature. For our experiments, we extracted a total of 15 different frequency features from pose-based time series, such as the horizontal displacement between the left hand and left shoulder, vertical displacement between the left hand and the left shoulder, and various angles, like the angle of the left elbow.

4 Classification: Frequency Forest

We adopt the forest-based method presented by O’Hara and Draper (2012) for our classification model. The data used, however, are frequency-domain features. Based on the algorithm presented, trees are constructed by incrementally adding training samples to an initially empty root node. Once a node becomes large enough and the splitting criteria are met, an element of the node is selected to be the pivot item and its distance to each item in the node is computed. Items having a distance to the pivot that is less than or equal to some threshold are added to the left child node, and the rest to the right. The now-empty node is marked as a splitting

node and all subsequently added instances are forwarded to one of its children after being compared to the pivot. The process recurses to form a tree. In the end, all interior nodes are splitting nodes and all leaf nodes contain neighboring samples.

In our experiments, we use the Entropy splitting method, which dictates that a tree’s node is only split when the distribution of distances between the items in the node falls below some empirically determined entropy threshold t_e (O’Hara and Draper 2012). (For our system, $t_e = 1.79$.) Furthermore, the Euclidean distance is used to measure the distance between two frequency features.

We build a frequency tree in the forest for each frequency feature. Each tree is designed to return the top five nearest neighbors for each test instance. However, only results from trees that have a distinct dominating label among the neighbors are considered for the final prediction. That is, we only consider results from trees where at least three of the top five neighbors share the same label. We believe this helps to prevent trees with weak correlations with the test instance from confusing the final prediction. Once all trees in the forest have voted, we pool all valid results together and return the most popular label as the final prediction.

5 Experiments

Recognition Experiment

For recognition, we tested on the KTH Actions dataset. The set contains six different actions (*boxing*, *handclapping*, *handwaving*, *jogging*, *running*, and *walking*), each performed several times by 25 different actors in four distinct scenarios: outdoors ($s1$), outdoors with scale variation ($s2$), outdoors with different clothes ($s3$), and indoors ($s4$) (Schuldt, Laptev, and Caputo 2004). There are a total of 2,391 action sequences.

We followed the same partitioning scheme outlined by Schuldt, Laptev, and Caputo (2004) to divide the dataset with respect to actors. The training and validation sets each contain video sequences for 8 unique actors. The test set contains sequences from the 9 remaining actors. However, since our method does not differentiate between the training and validation phase, both the training and validation sets (16 actors) were used to generate training examples.

Results from our recognition experiment can be seen in Table 1. Using frequency features derived from optical flow and articulated pose features, we achieved an overall accuracy rate of 82.7%. While this is not competitive with known state-of-the-art performances (see Table 2), it does show that our method works as it yields good performance rate.

Furthermore, the results show that we do well in recognizing the different hand-movement actions (e.g. *boxing*, *handclapping*, and *handwaving*). Our system also did a good job at recognizing *walking*.

Our biggest hurdle comes from distinguishing between *running* and *jogging*. Our system mistakenly labelled 38% of the run videos as jogging. This confusion can partly be explained by the fact that we were operating in the frequency domain, where the rate of body movements between running and jogging are often hard to differentiate. This is especially

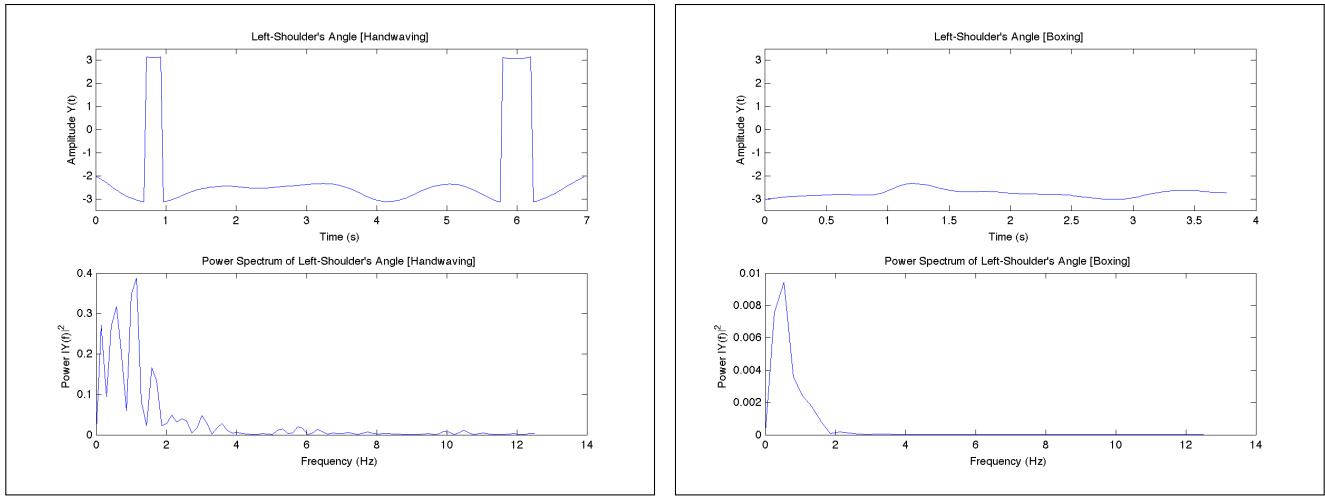


Figure 3: The smoothed left shoulder’s angle for two different action sequences (*handwaving* on the left and *boxing* on the right) and the corresponding power spectra in the frequency domain.

	bx	cl	wv	jg	rn	wk
box	0.81	0.12	0.06	0.00	0.00	0.01
clap	0.06	0.91	0.02	0.00	0.01	0.00
wave	0.02	0.13	0.85	0.00	0.00	0.00
jog	0.00	0.00	0.00	0.84	0.08	0.08
run	0.00	0.00	0.00	0.38	0.62	0.01
walk	0.01	0.00	0.00	0.04	0.01	0.94

Table 1: Performance confusion matrix for our frequency forest model on the KTH Actions dataset. The overall accuracy rate is 82.7%.

Method	Accuracy (%)
Schuld et al. [2004]	71.7%
Dollar et al. [2005]	81.2%
Frequency Forest	82.7%
O’hara and Draper [2012]	97.9%
Sadanand and Corso [2012]	98.2%

Table 2: Recognition accuracies on the KTH Actions dataset for various known methods, using Schuld’s training/testing partitioning of the data.

true across different actors, where one person’s rate of jogging may be indistinguishable from another person’s rate of running. Perhaps in cases like these, other methods might profitably augment frequency-based classification features.

Robustness to Variability

To show robustness, we set up an experiment that involves training and testing on different sets of data that vary in size and complexity/variability. These sets are formed by mixing and combining data from different video scenarios in the KTH Actions dataset.

Intuitively, we expect that if the test set is held constant,

then the increase in size and variability of the training set would result in more informative examples for the classifiers to train on. This should generally produce better and more reliable action models and positively influence recognition performance.

On the other hand, if we hold the training set constant and increase the size and variability of the test set, then it should make the test set more challenging for the classifiers. This should in turn negatively affect performance. However, if the features used to train the models are robust to variability, then one would expect the negative impact of an increase in variability in the test set to be small. That is to say recognition performance using these robust features should stay relatively stable and does not significantly decrease as the variability of the data in the test set increases.

For this experiment, we created different training sets by mixing data from different scenarios for all different combinations of scenarios. For testing, we adopted the four scenario configurations used by Schuld, Laptev, and Caputo (2004): $\{s1\}$, $\{s1, s4\}$, $\{s1, s3, s4\}$, $\{s1, s2, s3, s4\}$. Clearly, test configurations that include more scenarios are *harder* not only because of an increase in the number of test instances, but also due to the increase in variability between data from different scenarios. Furthermore, we retain the same actor-based partitioning scheme as above, which means we never train and test on the same actor. Thus, there is always an inherent variability in the actor involved, even for data from the same scenario.

Figure 5 shows the recognition accuracy for four representative training configurations. Complete results for all training combinations are shown in Table 3. Each accuracy rate given is the average of three independent runs.

Results shown in Figure 5 are representative of the general trends that we observed, and also validate our expectations of the system. As we increase the size and variability of the test set, performance tends to drop across the board, regardless of the training configuration. However, this de-

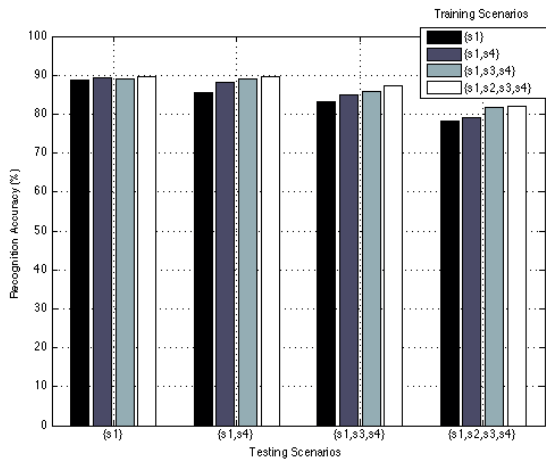


Figure 5: Recognition accuracies for four representative training configurations under different test sets. Performance for each training configuration drops slightly as the test set increases in size and complexity. Results shown are averaged across three independent runs

crease in performance is relatively small, which supports the robustness property of our features. In addition, the graph also shows that increasing the training set to cover more scenarios does positively increase performance slightly. That is to say, more training does help.

The complete set of results in Table 3 also generally follows the described trends except for a few odd cases, most of which involve the inclusion of scenario $s2$ in training and/or testing. This is not very surprising as scenario $s2$ is significantly different than the others. We think that the variability introduced by the variation in scales in $s2$ is often so overwhelming that it dominates, and perhaps alters, the actual underlying signal of the action, resulting in something that looks very much different than the same action from another scenario. As such, we believe that it would take, on average, more training data to successfully learn a good action model when using examples from $s2$. Hence, results from small training sets containing $s2$ may be unpredictable and difficult to interpret.

Additionally, there are few other odd cases that slightly contradict the general trends due to the overlapping, or lack thereof, between scenarios in the train and test sets. For example, when training on scenario $s4$ and testing on scenario $s1$, we performed at 85.0%. When we increase the test set to include another scenario, one would expect the performance to either remain stable or decrease. However, when the added scenario is also $s4$, generating a test set that includes $\{s1, s4\}$, the performance actually increased slightly to 86.4%, contradicting the trend. This, of course, can be explained by the overlapping of scenario $s4$ in both the train and test configurations.

Training Scenarios	Testing Scenarios			
	S1	S2	S3	S4
{s1}	88.7	85.4	83.3	78.2
{s2}	55.8	61.8	59.5	63.8
{s3}	88.8	85.2	83.0	76.8
{s4}	85.0	86.4	82.3	75.8
{s1, s2}	85.4	83.8	83.4	79.8
{s1, s3}	89.3	87.5	85.2	78.4
{s1, s4}	89.3	88.2	85.0	79.1
{s2, s3}	87.0	86.1	82.8	79.6
{s2, s4}	81.7	85.3	82.6	80.3
{s3, s4}	88.7	88.1	84.8	78.1
{s1, s2, s3}	88.8	86.9	84.9	80.6
{s1, s2, s4}	89.5	88.6	86.2	80.6
{s1, s3, s4}	89.9	88.3	86.8	79.1
{s2, s3, s4}	89.0	89.0	85.8	81.8
{s1, s2, s3, s4}	89.6	89.5	87.3	82.1

Table 3: Recognition accuracies for different combinations of train and test sets. For each row of training configuration, increasing the test set complexity and size does not significantly affect performance. Note: $\mathbf{S1} = \{s1\}$, $\mathbf{S2} = \{s1, s4\}$, $\mathbf{S3} = \{s1, s3, s4\}$, $\mathbf{S4} = \{s1, s2, s3, s4\}$. Results shown are averaged across three independent runs. Bold font indicates the best score in a given column.

6 Conclusion

This paper described a simple and efficient method to reduce the effects of data variability on recognition performance. We demonstrated how commonly used features like optical flow and articulated pose can be transformed into frequency-domain features that are less susceptible to variability in motions, viewpoints, dynamic backgrounds, and other challenging conditions. Furthermore, we showed how these frequency features can be used with a forest-based classifier to produce good and robust results on the KTH Actions dataset.

Although our method did not achieve the best performance in comparison with other known state-of-the-art systems (Table 2), we did achieve consistent performance that did not decrease much when variability in the test set increases. We attribute this stability in performance of our system to the robustness of our frequency-domain features.

We note that the presented method is not limited to the optical flow and articulated pose features that are described in the paper. We believe this technique can be used with any temporal data series, and plan to show the benefits of operating in the frequency domain for other recently developed action recognition features. In addition, we will further explore the extent to which frequency domain features are robust to variability. One of the first steps in this line of work will be to formally define the meaning of feature robustness and to formulate a way to accurately measure it. We are also interested in evaluating for robustness using cross-dataset testing, similar to the experiment described by Cao, Liu, and Huang (2010). We believe that our robust frequency features should do well in recognizing the same action across many different datasets and plan to test this hypothesis in future work.

References

- Aggarwal, J., and Cai, Q. 1999. Human Motion Analysis: A Review. *Computer Vision and Image Understanding* 73(3):428–440.
- Campbell, L., and Bobick, A. 1995. Recognition of human body motion using phase space constraints. In *Proceedings of IEEE International Conference on Computer Vision*, 624–630. IEEE Comput. Soc. Press.
- Cao, L.; Liu, Z.; and Huang, T. S. 2010. Cross-dataset action detection. In *2010 IEEE Computer Society Conference on Computer Vision and Pattern Recognition*, 1998–2005. IEEE.
- Danafar, S., and Gheissari, N. 2007. Action recognition for surveillance applications using optic flow and SVM. In *Proceedings of the 8th Asian conference on Computer vision - Volume Part II*, 457–466.
- Gavrila, D. 1999. The Visual Analysis of Human Movement: A Survey. *Computer Vision and Image Understanding* 73(1):82–98.
- Kerr, W.; Tran, A.; and Cohen, P. 2011. Activity Recognition with Finite State Machines. In *Twenty-Second International Joint Conference on Artificial Intelligence*, 1348–1353.
- Li, W.; Zhang, Z.; and Liu, Z. 2010. Action recognition based on a bag of 3D points. In *2010 IEEE Computer Society Conference on Computer Vision and Pattern Recognition - Workshops*, 9–14. IEEE.
- Liu, C. 2009. *Beyond pixels: exploring new representations and applications for motion analysis*. Phd dissertation, Massachusetts Institute of Technology.
- Lui, Y. M.; Beveridge, J. R.; and Kirby, M. 2010. Action classification on product manifolds. In *2010 IEEE Computer Society Conference on Computer Vision and Pattern Recognition*, 833–839. IEEE.
- Lv, F., and Nevatia, R. 2007. Single View Human Action Recognition using Key Pose Matching and Viterbi Path Searching. In *2007 IEEE Conference on Computer Vision and Pattern Recognition*, 1–8. IEEE.
- Moeslund, T. B., and Granum, E. 2001. A survey of computer vision-based human motion capture. *Computer Vision and Image Understanding* 81(3):231–268.
- O’Hara, S., and Draper, B. A. 2012. Scalable action recognition with a subspace forest. In *2012 IEEE Conference on Computer Vision and Pattern Recognition*, 1210–1217. IEEE.
- Ryoo, M., and Chen, C. 2010. An overview of contest on semantic description of human activities (SDHA) 2010. In *Proceedings of the 20th International conference on Recognizing patterns in signals, speech, images, and videos*, 270–285.
- Sadanand, S., and Corso, J. J. 2012. Action bank: A high-level representation of activity in video. In *2012 IEEE Conference on Computer Vision and Pattern Recognition*, 1234–1241. IEEE.
- Schuldt, C.; Laptev, I.; and Caputo, B. 2004. Recognizing human actions: a local SVM approach. In *Proceedings of the 17th International Conference on Pattern Recognition, 2004. ICPR 2004.*, 32–36 Vol.3. IEEE.
- Sheikh, Y.; Sheikh, M.; and Shah, M. 2005. Exploring the space of a human action. In *Tenth IEEE International Conference on Computer Vision (ICCV’05) Volume 1*, 144–149 Vol. 1. IEEE.
- Sung, J.; Ponce, C.; Selman, B.; and Saxena, A. 2011. Unstructured Human Activity Detection from RGBD Images. In *2012 IEEE International Conference on Robotics and Automation*.
- Turaga, P.; Chellappa, R.; Subrahmanian, V.; and Udrea, O. 2008. Machine Recognition of Human Activities: A Survey. *IEEE Transactions on Circuits and Systems for Video Technology* 18(11):1473–1488.
- Vondrick, C.; Ramanan, D.; and Patterson, D. 2010. Efficiently scaling up video annotation with crowdsourced marketplaces. In *Proceedings of the 11th European conference on Computer vision: Part IV*, 610–623. Heraklion, Crete, Greece: Springer-Verlag.
- Wang, J. J., and Singh, S. 2003. Video analysis of human dynamics a survey. *Real-Time Imaging* 9(5):321–346.
- Wang, J. 2012. Mining actionlet ensemble for action recognition with depth cameras. In *2012 IEEE Conference on Computer Vision and Pattern Recognition*, 1290–1297. IEEE.
- Xia, L.; Chen, C.; and Aggarwal, J. 2012. View Invariant Human Action Recognition Using Histograms of 3D Joints. In *IEEE CVPR Workshop on Human Activity Understanding from 3D Data*.
- Yacoob, Y., and Black, M. 1998. Parameterized modeling and recognition of activities. In *Sixth International Conference on Computer Vision (IEEE Cat. No.98CH36271)*, 120–127. Narosa Publishing House.
- Yang, Y., and Ramanan, D. 2011. Articulated pose estimation with flexible mixtures-of-parts. In *CVPR 2011*, 1385–1392. IEEE.