# Action Recognition Based on Joint Trajectory Maps Using Convolutional Neural Networks

**Pichao Wang[1*], Zhaoyang Li[2*], Yonghong Hou[2†] and Wanqing Li[1]**

[1]Advanced Multimedia Research Lab, University of Wollongong, Australia

[2]School of Electronic Information Engineering, Tianjin University, China

pw212@uowmail.edu.au,lizhaoyang@tju.edu.cn, houroy@tju.edu.cn, wanqing@uow.edu.au

## Abstract

Recently, Convolutional Neural Networks (ConvNets) have shown promising performances in many computer vision tasks, especially image-based recognition. How to effectively use ConvNets for video-based recognition is still an open problem. In this paper, we propose a compact, effective yet simple method to encode spatio-temporal information carried in $3D$ skeleton sequences into multiple $2D$ images, referred to as Joint Trajectory Maps (JTM), and ConvNets are adopted to exploit the discriminative features for real-time human action recognition. The proposed method has been evaluated on three public benchmarks, i.e., MSRC-12 Kinect gesture dataset (MSRC-12), G3D dataset and UTD multimodal human action dataset (UTD-MHAD) and achieved the state-of-the-art results.

## 1 Introduction

Recognition of human actions from RGB-D (Red, Green, Blue and Depth) data has attracted increasing attention in multimedia signal processing in recent years due to the advantages of depth information over conventional RGB video, e.g. being insensitive to illumination changes. Since the first work of such a type [9] reported in 2010, many methods [17; 12; 23; 10] have been proposed based on specific hand-crafted feature descriptors extracted from depth. With the recent development of deep learning, a few methods [18; 19] have been developed based on Convolutional Neural Networks (ConvNets). A common and intuitive method to represent human motion is to use a sequence of skeletons. With the development of the cost-effective depth cameras and algorithms for real-time pose estimation [14], skeleton extraction has become more robust and many hand-designed skeleton features [22; 24; 5; 20; 16] for action recognition have been proposed. Recently, Recurrent Neural Networks (RNNs) [3; 15; 28; 13] have also been adopted for action recognition from skeleton data. The hand-crafted features are always shallow and dataset-dependent. RNNs tend to overemphasize

the temporal information especially when the training data is not sufficient, leading to overfitting. In this paper, we present a compact, effective yet simple method that encodes the joint trajectories into texture images, referred to as Joint Trajectory Maps (JTM), as the input of ConvNets for action recognition. In this way, the capability of the ConvNets in learning discriminative features can be fully exploited [25].

One of the challenges in action recognition is how to properly model and use the spatio-temporal information. The commonly used bag-of-words model tends to overemphasize the spatial information. On the other hand, Hidden Markov Model (HMM) or RNN based methods are likely to overstress the temporal information. The proposed method addresses this challenge in a different way by encoding as much the spatio-temporal information as possible (without a need to decide which one is important and how important it is) into images and letting the CNNs to learn the discriminative one. This is the key reason that the proposed method outperformed previous ones. In addition, the proposed encoding method can be extended to online recognition due to the accumulative nature of the encoding process. Furthermore, such encoding of spatio-temporal information into images allows us to leverage the advanced methods developed for image recognition.

## 2 The Proposed Method

The proposed method consists of two major components, as illustrated in Fig. 1, three ConvNets and the construction of three JTMs as the input of the ConvNets in three orthogonal planes from the skeleton sequences. Final classification of a given test skeleton sequence is obtained through a late fusion of the three ConvNets. The main contribution of this paper is on the construction of suitable JTMs for the ConvNets to learn discriminative features.

We argue that an effective JTM should have the following properties to keep sufficient spatial-temporal information of an action:

- The joints or group of joints should be distinct in the JTM such that the spatial information of the joints is well reserved.

- The JTM should encode effectively the temporal evolution, i.e. trajectories of the joints, including the direction and speed of joint motions.

---

[*]Both authors contributed equally to this work
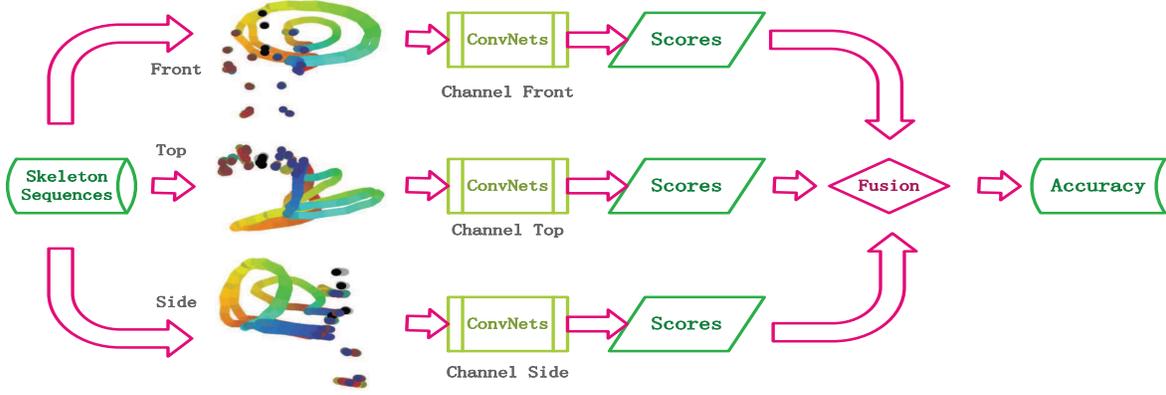
[†]Corresponding author

Figure 1: The framework of the proposed method.

- The JTM should be able to encode the difference in motion among the different joints or parts of the body to reflect how the joints are synchronized during the action.

Specifically, JTM can be recursively defined as follows

$$JTM_i = JTM_{i-1} + f(i) \qquad (1)$$

where $f(i)$ is a function encoding the spatial-temporal information at frame or time-stamp $i$. Since JTM is accumulated over the period of an action, $f(i)$ has to be carefully defined such that the JTM for an action sample has the required properties and the accumulation over time has little adverse impact on the spatial-temporal information encoded in the JTM. We propose in this paper to use hue, saturation and brightness to encode the spatial-temporal motion patterns.

### 2.1  Joint Trajectory Maps

Assume an action $H$ has $n$ frames of skeletons and each skeleton consists of $m$ joints. The skeleton sequence is denoted as $H = \{F_1, F_2, ..., F_n\}$, where $F_i = \{P_1^i, P_2^i, ..., P_m^i\}$ is a vector of the joint coordinates at frame $i$, and $P_j^i$ is the $3D$ coordinates of the $j$th joint in frame $i$. The skeleton trajectory $T$ for an action of $n$ frames consists of the trajectories of all joints and is defined as:

$$T = \{T_1, T_2, \cdots, T_i, \cdots, T_{n-1}\} \qquad (2)$$

where $T_i = \{t_1^i, t_2^i, ..., t_m^i\} = F_{i+1} - F_i$ and the $k$th joint trajectory is $t_k^i = P_k^{i+1} - P_k^i$. At this stage, the function $f(i)$ is the same as $T_i$, that is,

$$f(i) = T_i = \{t_1^i, t_2^i, ..., t_m^i\}. \qquad (3)$$

The skeleton trajectory is projected to the three orthogonal planes, i.e. three Cartesian planes, to form three JTMs. Fig. 2 shows the three projected trajectories of the right hand joint for action "right hand draw circle (clockwise)" in the UTD-MHAD dataset. From these JTMs, it can be seen that the spatial information of this joint is preserved but the direction of the motion is lost.
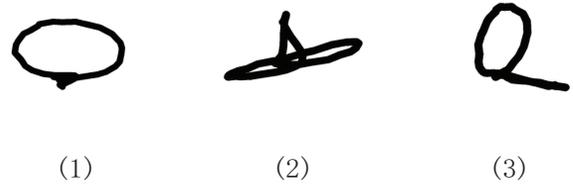


Figure 2: The trajectories projected onto three Cartesian planes for action "right hand draw circle (clockwise)" in UTD-MHAD [2]: (1) the front plane; (2) the top plane; (3) the side plane.

### 2.2  Encoding Joint Motion Direction

To capture the motion information in the JTM, it is proposed to use hue to represent the motion direction. Different kinds of colormaps can be chosen. In this paper, the jet colormap, ranging from blue to red, and passing through the colors cyan, yellow, and orange, was adopted. Assume the color of a joint trajectory is $C$ and the length of the trajectory $L$, and let $C_l, l \in (0, L)$ be the color at position $l$. For the $q^{th}$ trajectory $T_q$ from 1 to $n-1$, a color $C_l$, where $l = \frac{q}{n-1} \times L$ is specified to the joint trajectory, making different trajectories have their own color corresponding to their temporal positions in the sequence as illustrated in Fig. 3. Herein, the trajectory with color is denoted as $C\_t_k^i$ and the function $f(i)$ is updated to:

$$f(i) = \{C\_t_1^i, C\_t_2^i, ..., C\_t_m^i\}. \qquad (4)$$

This ensures that different actions are encoded to a same length colormap. The effects can be seen in Fig. 4, subfigures (1) to (2). Even though the same actions with different number of cycles will be encoded into different color shapes, the direction can still be reflected in color variation and the differences between actions can still be captured due to the different spatial information.
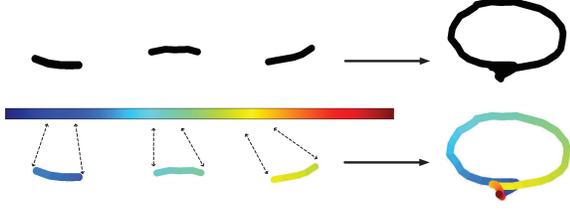
Figure 3: The trajectories of different body parts have their different colors reflecting the temporal orders.

## 2.3 Encoding Body Parts

To distinguish different body parts, multiple colormaps are employed. There are many ways to achieve this. For example, each joint is assigned to one colormap, or several groups of joints are assigned to different colormaps randomly. Considering arms and legs often have more motion than other body parts, we empirically generate three colormaps ($C1, C2, C3$) to encode three body parts. $C1$ is used for the left body part (consisting of left shoulder, left elbow, left wrist, left hand, left hip, left knee, left ankle and left foot), $C2$ for the right body part ( consisting of right shoulder, right elbow, right wrist, right hand, right hip, right knee, right ankle and right foot), and $C3$ for the middle body part (consisting of head, neck, torso and hip center). $C1$ is the same as $C$, i.e. the jet colormap, $C2$ is a reversed colormap of $C1$, and $C3$ is a colormap ranging from light gray to black. Here, the trajectory encoded by multiple colormaps is denoted as $MC\_t_k^i$, and the function $f(i)$ is formulated as:

$$f(i) = \{MC\_t_1^i, MC\_t_2^i, ..., MC\_t_m^i\}. \tag{5}$$

The effects can be seen in Fig. 4, sub-figures (2) to (3).

## 2.4 Encoding Motion Magnitude

Motion magnitude is one of the most important factors in human motion. For one action, large magnitude of motion usually indicates more motion information. In this paper, it is proposed to encode the motion magnitude of joints into the saturation and brightness components, so that such encoding not only encodes the motion but also enriches the texture of trajectories which are expected to be beneficial for ConvNets to learn discriminative features. For joints with high motion magnitude or speed, high saturation will be assigned as high motion usually carries more discriminative information. Specifically, the saturation is set to range from $s_{min}$ to $s_{max}$. Given a trajectory, its saturation $S_j^i$ in $HSV$ color space could be calculated as

$$S_j^i = \frac{v_j^i}{max\{v\}} \times (s_{max} - s_{min}) + s_{min} \tag{6}$$

where $v_j^i$ is the $j$th joint speed at the $i$th frame.

$$v_j^i = \|P_j^{i+1} - P_j^i\|_2 \tag{7}$$

The trajectory adjusted by saturation is denoted as $MC_s\_t_k^i$ and the function $f(i)$ is refined as:

$$f(i) = \{MC_s\_t_1^i, MC_s\_t_2^i, ..., MC_s\_t_m^i\} \tag{8}$$

The encoding effect can be seen in Figure 4, sub-figures (3) to (4), where the slow motion becomes diluted (e.g. trajectory of knees and ankles) while the fast motion becomes saturated (e.g. the green part of the circle).

To further enhance the motion patterns in the JTM, the brightness is modulated by the speed of joints so that motion information is enhance in the JTM by rapidly changing the brightness according to the joint speed. In particular, the brightness is set to range from $b_{min}$ to $b_{max}$. Given a trajectory $t_j^i$ whose speed is $v_j^i$, its brightness $B_j^i$ in the $HSV$ color space is calculated as

$$B_j^i = \frac{v_j^i}{max\{v\}} \times (b_{max} - b_{min}) + b_{min} \tag{9}$$

The trajectory adjusted by brightness is denoted as $MC_b\_t_k^i$ and the function $f(i)$ is updated to:

$$f(i) = \{MC_b\_t_1^i, MC_b\_t_2^i, ..., MC_b\_t_m^i\}. \tag{10}$$

The effect can be seen in Fig 4, sub-figures (3) to (5), where texture becomes apparent (e.g. the yellow parts of the circle). Finally, motion magnitude is encoded with saturation and brightness together. The trajectory is denoted as $MC_{sb}\_t_k^i$ and the function $f(i)$ is refined as:

$$f(i) = \{MC_{sb}\_t_1^i, MC_{sb}\_t_2^i, ..., MC_{sb}\_t_m^i\}. \tag{11}$$

As illustrated in Fig. 4, sub-figures(3) to (6), it not only enriches the texture information but also highlights the faster motion.
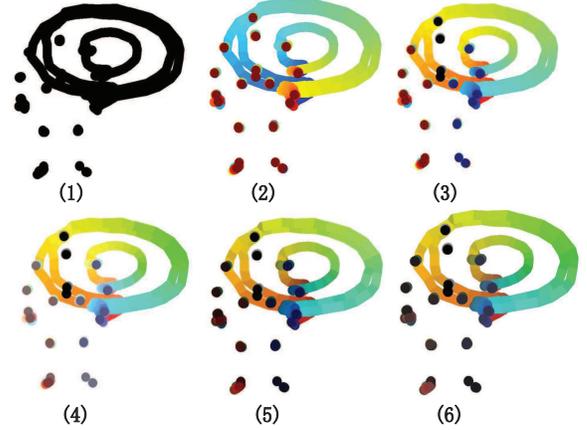


Figure 4: Illustration of visual differences for different techniques in JTM.

## 2.5 Training and Recognition

In the experiments, the layer configuration of the three ConvNets was same as the one in [8]. The implementation was

derived from the publicly available Caffe toolbox [7] based on one NVIDIA GeForce GTX TITAN X card and the pre-trained models over ImageNet [8] were used for initialization in training. The network weights are learned using the mini-batch stochastic gradient descent with the momentum being set to 0.9 and weight decay being set to 0.0005. At each iteration, a mini-batch of 256 samples is constructed by sampling 256 shuffled training JTMs. All JTMs are resized to $256 \times 256$. The learning rate is to $10^{-3}$ for fine-tuning and then it is decreased according to a fixed schedule, which is kept the same for all training sets. For each ConvNet the training undergoes 100 cycles and the learning rate decreases every 20 cycles. For all experiments, the dropout regularisation ratio was set to 0.5 in order to reduce complex co-adaptations of neurons in nets. Three ConvNets are trained on the JTMs in the three Cartesian planes and the final score for a test sample are the averages of the outputs from the three ConvNets. The testing process can easily achieved real-time speed (average 0.36 seconds/sample).

## 3 Experimental Results

The proposed method was evaluated on three public benchmark datasets: MSRC-12 Kinect Gesture Dataset [4], G3D [1] and UTD-MHAD [2]. Experiments were conducted to evaluate the effectiveness of each encoding scheme in the proposed method and the final results were compared with the state-of-the-art reported on the same datasets. In all experiments, the saturation and brightness covers the full range (from $0\% \sim 100\%$ mapped to $0 \sim 255$) in HSV color space.

### 3.1 Evaluation of Different Encoding Schemes

The effectiveness of different encoding schemes (corresponding to the sub-figures in 4) was evaluated on the G3D dataset using the front JTM and the recognition accuracies are listed in Table 1.

| Techniques | Accuracy (%) |
|---|---|
| Trajectory: $t_1^i$ | 63.64% |
| Trajectory: $C\_t_1^i$ | 74.24% |
| Trajectory: $MC\_t_1^i$ | 78.48% |
| Trajectory: $MC_s\_t_1^i$ | 81.82% |
| Trajectory: $MC_b\_t_1^i$ | 82.12% |
| Trajectory: $MC_{sb}\_t_1^i$ | 85.45% |

Table 1: Comparisons of the different encoding schemes on the G3D dataset using the JTM projected to the front plane alone.

From this Table it can be seen that the proposed encoding techniques effectively captures the spatio-temporal information and the ConvNets are able to learn the discriminative features from the JTM for action recognition.

### 3.2 MSRC-12 Kinect Gesture Dataset

MSRC-12 [4] is a relatively large dataset for gesture/action recognition from 3D skeleton data captured by a Kinect sensor. The dataset has 594 sequences, containing 12 gestures by 30 subjects, 6244 gesture instances in total. The 12 gestures

are: "lift outstretched arms", "duck", "push right", "goggles", "wind it up", "shoot", "bow", "throw", "had enough", "beat both", "change weapon" and "kick". For this dataset, cross-subjects protocol is adopted, that is odd subjects for training and even subjects for testing. Table 2 lists the performance of the proposed method and the results reported before.

| Method | Accuracy (%) |
|---|---|
| HGM [21] | 66.25% |
| ELC-KSVD [27] | 90.22% |
| Cov3DJ [6] | 91.70% |
| Proposed Method | **93.12%** |

Table 2: Comparison of the proposed method with the existing methods on the MSRC-12 Kinect gesture dataset.

The confusion matrix is shown in figure 5. From the confusion matrix we can see that the proposed method distinguishes most of actions very well, but it is not very effective to distinguish "goggles" and "had enough" which shares the similar appearance of JTM probably caused by 3D to 2D projection.
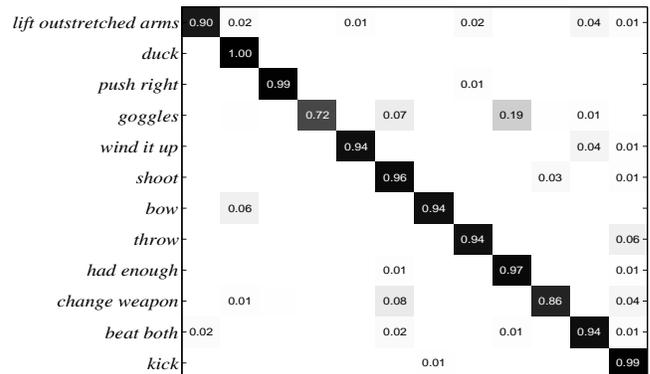


Figure 5: The confusion matrix of the proposed method for MSRC-12 Kinect gesture dataset.

### 3.3 G3D Dataset

Gaming 3D Dataset (G3D) [1] focuses on real-time action recognition in gaming scenario. It contains 10 subjects performing 20 gaming actions: "punch right", "punch left", "kick right", "kick left", "defend", "golf swing", "tennis swing forehand", "tennis swing backhand", "tennis serve", "throw bowling ball", "aim and fire gun", "walk", "run", "jump", "climb", "crouch", "steer a car", "wave", "flap" and "clap". For this dataset, the first 4 subjects were used for training, the fifth for validation and the remaining 5 subjects for testing as configured in [11].

Table 3 compared the performance of the proposed method and that reported in [11].

The confusion matrix is shown in figure 6. From the confusion matrix we can see that the proposed method recognizes most of actions well. Compared with LRBM, our proposed method outperforms LRBM in spatial information mining.

| Method | Accuracy (%) |
|---|---|
| LRBM [11] | 90.50% |
| Proposed Method | **94.24%** |

Table 3: Comparison of the proposed method with previous methods on G3D Dataset.

LRBM confused mostly the actions between "tennis swing forehand" and "bowling", "golf" and "tennis swing backhand", "aim and fire gun" and "wave", "jump" and "walk", however, these actions were quite well distinguished in our method because of the good spatial information exploitation in our method. As for "aim and fire gun" and "wave", our method could not distinguish them well before encoding the motion magnitude, which means the temporal information enhancement procedure is effective. However, in our method, "tennis swing forehand" and "tennis swing backhand" are confused. It's probably because the front and side projections of body shape of the two actions are too similar, and scores fusion is not very effective to improve each other.
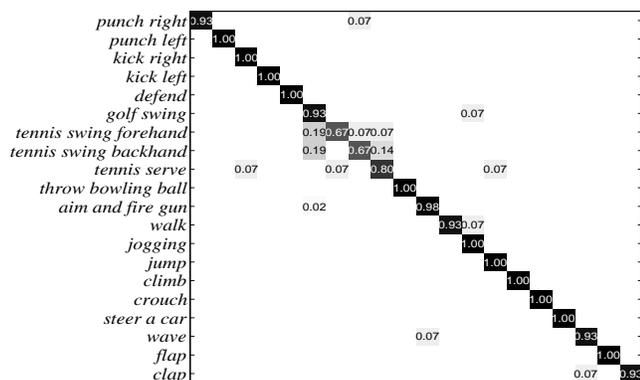


Figure 6: The confusion matrix of the proposed method for G3D Dataset.

## 3.4 UTD-MHAD

UTD-MHAD [2] is one multimodal action dataset, captured by one Microsoft Kinect camera and one wearable inertial sensor. This dataset contains 27 actions performed by 8 subjects (4 females and 4 males) with each subject perform each action 4 times. After removing three corrupted sequences, the dataset includes 861 sequences. The actions are: "right arm swipe to the left", "right arm swipe to the right", "right hand wave", "two hand front clap", "right arm throw", "cross arms in the chest", "basketball shoot", "right hand draw x", "right hand draw circle (clockwise)", "right hand draw circle (counter clockwise)", "draw triangle", "bowling (right hand)", "front boxing", "baseball swing from right", "tennis right hand forehand swing", "arm curl (two arms)", "tennis serve", "two hand push", "right hand know on door", "right hand catch an object", "right hand pick up and throw", "jogging in place", "walking in place", "sit to stand", "stand to sit", "forward lunge (left foot forward)" and "squat (two arms stretch out)". It covers sport actions (e.g. "bowling", "tennis serve" and "baseball swing"), hand gestures (e.g. "draw

X", "draw triangle", and "draw circle"), daily activities (e.g. "knock on door", "sit to stand" and "stand to sit") and training exercises (e.g. "arm curl", "lung" and "squat"). For this dataset, cross-subjects protocol is adopted as in [2], namely, the data from the subject numbers 1, 3, 5, 7 used for training while 2, 4, 6, 8 used for testing.

Table 4 compared the performance of the proposed method and that reported in [2].

| Method | Accuracy (%) |
|---|---|
| Kinect & Inertial [2] | 79.10% |
| Proposed Method | **85.81%** |

Table 4: Comparison of the proposed method with previous methods on UTD-MHAD Dataset.

Please notice that the method used in [2] is based on Depth and Inertial sensor data, not skeleton data alone.
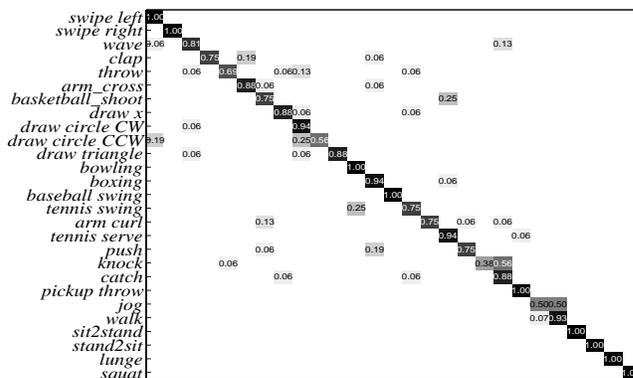


Figure 7: The confusion matrix of the proposed method for UTD-MHAD.

The confusion matrix is shown in figure 7. This dataset is much more challenging compared to previous two datasets. From the confusion matrix we can see that the proposed method can not distinguish some actions well, for example, "jog" and "walk". A probable reason is that the proposed encoding process is also a normalization process along temporal axis (Section 3.2). The actions "jog" and "walk" will be normalized to have a very similar JTM after the encoding.

## 4 Conclusion

This paper addressed the problem of human action recognition by applying ConvNets to skeleton sequences. We proposed an effective method to encode the joints trajectories to JTM where the motion information can be encoded into texture patterns. ConvNets learn discriminative features from these maps for real-time human action recognition. The experimental results showed that the techniques for encoding worked effectively. The proposed method can benefit from effective data augmentation process which would be our future work.

# 5  Acknowledgments

## References

[1] V. Bloom, D. Makris, and V. Argyriou. G3D: A gaming action dataset and real time action recognition evaluation framework. In *Proc. IEEE Computer Society Conference on Computer Vision and Pattern Recognition Workshops (CVPRW)*, pages 7–12, 2012.

[2] C. Chen, R. Jafari, and N. Kehtarnavaz. Utd-mhad: A multimodal dataset for human action recognition utilizing a depth camera and a wearable inertial sensor. In *Image Processing (ICIP), 2015 IEEE International Conference on*, pages 168–172, 2015.

[3] Y. Du, W. Wang, and L. Wang. Hierarchical recurrent neural network for skeleton based action recognition. In *Proc. IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 1110–1118, 2015.

[4] S. Fothergill, H. M. Mentis, S. Nowozin, and P. Kohli. Instructing people for training gestural interactive systems. In *ACM Conference on Computer-Human Interaction (ACM HCI)*, 2012.

[5] M. A. Gowayyed, M. Torki, M. E. Hussein, and M. El-Saban. Histogram of oriented displacements (HOD): Describing trajectories of human joints for action recognition. In *Proc. International Joint Conference on Artificial Intelligence (IJCAI)*, pages 1351–1357, 2013.

[6] M. E. Hussein, M. Torki, M. A. Gowayyed, and M. El-Saban. Human action recognition using a temporal hierarchy of covariance descriptors on 3D joint locations. In *Proc. International Joint Conference on Artificial Intelligence (IJCAI)*, pages 2466–2472, 2013.

[7] Y. Jia, E. Shelhamer, J. Donahue, S. Karayev, J. Long, R. B. Girshick, S. Guadarrama, and T. Darrell. Caffe: Convolutional architecture for fast feature embedding. In *Proc. ACM international conference on Multimedia (ACM MM)*, pages 675–678, 2014.

[8] A. Krizhevsky, I. Sutskever, and G. E. Hinton. Imagenet classification with deep convolutional neural networks. In *Proc. Annual Conference on Neural Information Processing Systems (NIPS)*, pages 1106–1114, 2012.

[9] W. Li, Z. Zhang, and Z. Liu. Action recognition based on a bag of 3D points. In *Proc. IEEE Computer Society Conference on Computer Vision and Pattern Recognition Workshops (CVPRW)*, pages 9–14, 2010.

[10] C. Lu, J. Jia, and C.-K. Tang. Range-sample depth feature for action recognition. In *Proc. IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 772–779, 2014.

[11] S. Nie, Z. Wang, and Q. Ji. A generative restricted boltzmann machine based method for high-dimensional motion data modeling. *Computer Vision and Image Understanding*, pages 14–22, 2015.

[12] O. Oreifej and Z. Liu. HON4D: Histogram of oriented 4D normals for activity recognition from depth sequences. In *Proc. IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 716–723, 2013.

[13] A. Shahroudy, J. Liu, T.-T. Ng, and G. Wang. NTU RGB+ D: A large scale dataset for 3D human activity analysis. In *Proc. IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2016.

[14] J. Shotton, A. Fitzgibbon, M. Cook, T. Sharp, M. Finocchio, R. Moore, A. Kipman, and A. Blake. Real-time human pose recognition in parts from single depth images. In *Proc. IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 1297–1304, 2011.

[15] V. Veeriah, N. Zhuang, and G.-J. Qi. Differential recurrent neural networks for action recognition. In *Proc. IEEE International Conference on Computer Vision (ICCV)*, pages 4041–4049, 2015.

[16] R. Vemulapalli, F. Arrate, and R. Chellappa. Human action recognition by representing 3D skeletons as points in a lie group. In *Proc. IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 588–595, 2014.

[17] J. Wang, Z. Liu, Y. Wu, and J. Yuan. Mining actionlet ensemble for action recognition with depth cameras. In *Proc. IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 1290–1297, 2012.

[18] P. Wang, W. Li, Z. Gao, C. Tang, J. Zhang, and P. O. Ogunbona. Convnets-based action recognition from depth maps through virtual cameras and pseudocoloring. In *Proc. ACM international conference on Multimedia (ACM MM)*, pages 1119–1122, 2015.

[19] P. Wang, W. Li, Z. Gao, J. Zhang, C. Tang, and P. Ogunbona. Action recognition from depth maps using deep convolutional neural networks. *Human-Machine Systems, IEEE Transactions on*, PP(99):1–12, 2015.

[20] P. Wang, W. Li, P. Ogunbona, Z. Gao, and H. Zhang. Mining mid-level features for action recognition based on effective skeleton representation. In *Proc. International Conference on Digital Image Computing: Techniques and Applications (DICTA)*, pages 1–8, 2014.

[21] S. Yang, C. Yuan, W. Hu, and X. Ding. A hierarchical model based on latent dirichlet allocation for action recognition. In *Pattern Recognition (ICPR), 2014 22nd International Conference on*, pages 2613–2618. IEEE, 2014.

[22] X. Yang and Y. Tian. Eigenjoints-based action recognition using Naive-Bayes-Nearest-Neighbor. In *Proc. International Workshop on Human Activity Understanding from 3D Data (HAU3D) (CVPRW)*, pages 14–19, 2012.

[23] X. Yang and Y. Tian. Super normal vector for activity recognition using depth sequences. In *Proc. IEEE International Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 804–811, 2014.

[24] M. Zanfir, M. Leordeanu, and C. Sminchisescu. The moving pose: An efficient 3D kinematics descriptor for low-latency action recognition and detection. In *Proc. IEEE International Conference on Computer Vision (ICCV)*, pages 2752–2759, 2013.

[25] M. D. Zeiler and R. Fergus. Visualizing and understanding convolutional networks. In *Proc. European Conference on Computer Vision (ECCV)*, pages 818–833. 2014.

[26] J. Zhang, W. Li, P. O. Ogunbona, P. Wang, and C. Tang. Rgb-d-based action recognition datasets: A survey. *Pattern Recognition*, 60:86–105, 2016.

[27] L. Zhou, W. Li, Y. Zhang, P. Ogunbona, D. T. Nguyen, and H. Zhang. Discriminative key pose extraction using extended lc-ksvd for action recognition. In *Proc. International Conference on Digital Image Computing: Techniques and Applications (DICTA)*, pages 1–8. IEEE, 2014.

[28] W. Zhu, C. Lan, J. Xing, W. Zeng, Y. Li, L. Shen, and X. Xie. Co-occurrence feature learning for skeleton based action recognition using regularized deep LSTM networks. In *The 30th AAAI Conference on Artificial Intelligence (AAAI)*, 2016.