

INVESTIGATION OF DIFFERENT SKELETON FEATURES FOR CNN-BASED 3D ACTION RECOGNITION

Zewei Ding, Pichao Wang*, Philip O. Ogunbona, Wanqing Li

Advanced Multimedia Research Lab, University of Wollongong, Australia
 {zd027, pw212, philipo, wanqing}@uow.edu.au

ABSTRACT

Deep learning techniques are being used in skeleton based action recognition tasks and outstanding performance has been reported. Compared with RNN based methods which tend to overemphasize temporal information, CNN-based approaches can jointly capture spatio-temporal information from texture color images encoded from skeleton sequences. There are several skeleton-based features that have proven effective in RNN-based and handcrafted-feature-based methods. However, it remains unknown whether they are suitable for CNN-based approaches. This paper proposes to encode five spatial skeleton features into images with different encoding methods. In addition, the performance implication of different joints used for feature extraction is studied. The proposed method achieved state-of-the-art performance on NTU RGB+D dataset for 3D human action analysis. An accuracy of 75.32% was achieved in Large Scale 3D Human Activity Analysis Challenge in Depth Videos.

Index Terms— Skeleton, 3D Action recognition, Convolutional Neural Networks

1. INTRODUCTION

Recognition of human actions has recently attracted increased interest because of its applicability in systems such as human-computer interaction, game control, and intelligent surveillance. With the development of cost-effective sensors such as Microsoft Kinect cameras, RGB-D-based recognition has almost become commonplace [1, 2, 3, 4]. Among the three most common input streams (RGB, depth, and skeleton), RGB is the most popular and widely studied. However, it suffers the challenge of pose ambiguity due to the loss of 3D information. On the other hand, depth and skeleton which capture 3D information of human bodies inherently overcome this challenge.

Skeleton has the advantage of being invariant to viewpoints or appearances compared with depth, thus suffering less intra-class variance [5]. Furthermore, learning over skeleton is simple because they are higher-level information

based on advanced pose estimation. The foregoing observations motivated the study of skeleton-based human action recognition in this paper.

The methods based on handcrafted skeleton features [6, 7, 8] have the drawback of dataset dependency while methods based on deep learning techniques have achieved outstanding performance. Currently, there are mainly two ways of using deep learning techniques to capture the spatio-temporal information in skeleton sequences; Recurrent Neural Networks (RNNs) and Convolutional Neural Networks (CNNs). RNNs are adopted to capture temporal information from extracted spatial skeleton features. The performance relies much on the effectiveness of the extracted spatial skeleton features due to the sequential flow of information. Moreover, the temporal information can be easily overemphasized especially when the training data is insufficient, leading to overfitting [9].

In contrast, CNNs directly extract information from texture images which are encoded from skeleton sequences. Wang et al [9] used Joint Trajectory Maps (JTM) to encode body joint trajectories (positions, motion directions, and motion magnitudes) of each time instance into HSV images. In the images, spatial information is represented by positions and the dynamics is represented by colors. Hou et al [10] adopted Skeleton Optical Spectra (SOS) to encode dynamic spatio-temporal information. Li et al [11] adopted joint distances as spatial features and a colorbar was used for color-encoding. In the images, textures of rows capture spatial information and textures of columns capture temporal information. Currently, the spatial features used for encoding are relatively simple (joints positions and pair-wise distances).

Following the CNN-based approach, this paper investigates encoding richer spatial features into texture color images, including features between two or more joints. Specifically, inspired by the work from Zhang et al [5], the encoding of the following five types of spatial features is studied: joint-joint distances (JJd), joint-joint orientations (JJo), joint-joint vectors (JJv), joint-line distances (JLd), line-line angles (LLa). Each kind of feature is encoded into images in two or more ways to further explore the spatio-temporal information. CNN is adopted to train and recognize corresponding actions and score fusion is used to make a final classification.

*Corresponding author

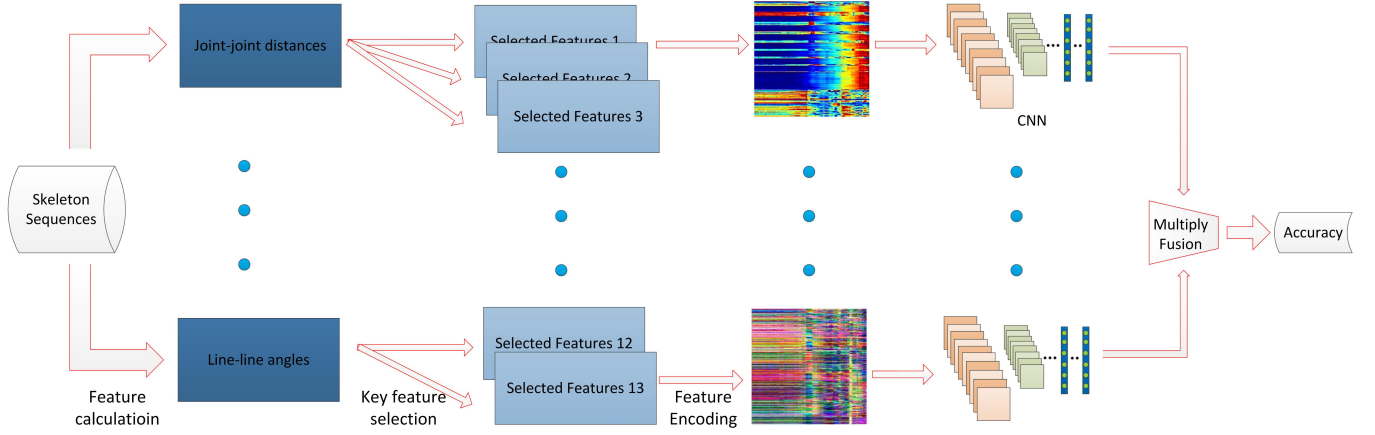


Fig. 1. The framework of the proposed method

The effectiveness of this kind of approach has been verified in [9, 11, 2]. The investigation is conducted on NTU RGB+D Dataset [12] and achieves state-of-the-art performance.

The rest of the paper is organized as follows. Section 2 introduces the proposed method and, in Section 3, experimental results and discussions are described. The conclusion and future work are presented in Section 4.

2. PROPOSED METHOD

As illustrated in Fig. 1, the proposed method consists of five main components, namely spatial feature extraction from input skeleton sequences, key feature selection, texture color image encoding from key features, CNN model training based on images, and the score fusion. There are five types of features extracted from all joint combinations including JJd, JJo, JJv, JLd and LLa. Key features of certain joint combinations are then chosen for color encoding. For each type of key features, there are multiple selection methods and encoding methods, resulting in a total of 13 types of images. CNN is trained on each kind of image, and the output scores of CNNs are fused into the final score for final recognition.

2.1. Feature extraction

The spatial features studied in this paper include joint-joint distances, joint-joint vectors, joint-joint orientations, joint-line distances and line-line angles which were introduced in [5]. In this paper, every action is assumed to be performed by two subjects, the main subject and an auxiliary other. In cases where there is only one person in the sequence, a 'shadow subject' copied from main subject is adopted. Suppose each subject has n joints, then in each frame there will be $N = 2 \times n$ joints. Let $p_j = (x, y, z)$, $j \in N$ denote the 3D coordinate (Euclidean space) of the j_{th} joint in a frame. The five features at frame t are calculated as follows:

$$JJd_{jk}^t = \|p_j^t - p_k^t\| \quad (1)$$

$$JJv_{jk}^t = p_j^t - p_k^t \quad (2)$$

$$JJo_{jk}^t = JJv_{jk}^t / JJd_{jk}^t \quad (3)$$

$$JLd_{jkm}^t = JJv_{jk}^t \otimes JJv_{jm}^t / JJd_{jk}^t \quad (4)$$

$$LLa_{jkmn}^t = \arccos(JJo_{jk}^t \odot JJo_{mn}^t) \quad (5)$$

where $j, k, m, n \in N$ are the joint indices, \otimes is cross product and \odot is dot product. Meanwhile, $j \neq k$ in equations (1-4), $j \neq m \neq k$ in equation (4), and $(i, k) \neq (m, n)$ in equation (5).

In total, there are $C_{50}^2 = 1225$ dimensions of the JJd feature, $3 \times 1225 = 3675$ dimensions of JJv and JJo features. There are also 1225 lines, resulting in $1225 \times 48 = 58800$ dimensions of JLd feature and $C_{1225}^2 = 749700$ dimensions of LLa feature. The resulting high dimensional feature space is neither cost-effective nor robust.

2.2. Feature selection

Feature selection is conducted by selecting key joints and key lines to reduce the number of combinations. The selection follows the principle that selected features should contain as much information as possible and be invariant to viewpoints and actions. Based on the observation that the motions are mainly located on the ends of skeletons and are usually locally sensitive, three strategies are proposed to select key joints for joint-joint feature calculation.

Joint strategy one (JS1): only the relations of joints within the same subject are considered, resulting in $2 \times C_{25}^2 = 600$ dimensional JJd feature. JS2: twelve joints from each subject are used, resulting in $C_{24}^2 = 276$ dimensional JJd feature. The joints start from 'middle of the spine' and are all two-steps away from the others. JS3: eleven joints from each subject are used, resulting in $C_{22}^2 = 231$ dimensional JJd feature.

The joints start from 'base of the spine' and are all two-steps away from the others.

Two strategies are used to select key lines. Line strategy one (LS1): adopting the method in [5] to select 39 lines from the main subject, resulting in 897 dimensional JLD feature and 741 dimensional LLa feature. LS2: using joints selected via JS3 to generate lines, and for each line the joints within two-step distance from end joints are used to calculate JLD feature, resulting in 570 dimensional JLD feature.

2.3. Color encoding

Inspired by [11], color images are used to encode the spatial features to capture temporal information. Specifically, each column in the image represents spatial features in a frame, and each row represents the sequence of a specific feature. In this way, the textures represent the spatio-temporal information of skeleton sequences. Given a skeleton sequence with T frames, N -dimensional features (scalar/vector) are extracted for each frame. The following three methods are used to encode the $N \times T$ feature into a $H(\text{height}) \times W(\text{width})$ sized color image (256×256 in this paper).

Encoding method one (EM1): for scalar features including JJD, JLD and LLa, the jet colorbar [11] is adopted to encode RGB channels jointly. The RGB value of pixel at h_{th} row and w_{th} column is

$$\overline{RGB}(h, w) = \text{colorbar}((f_h^w - \min F_h) / (\max(F_h) - \min F_h)) \quad (6)$$

where f_h^w is the value of the h_{th} feature at w_{th} frame, and $f_h^w = f_{N \times h/H}^{T \times w/W}$, i.e. the features are resized to $H * W$ using bilinear interpolation. $F_h = \{f_h^1, f_h^2, \dots, f_h^T\}$, $\text{colorbar}()$ is a mapping function which maps $[0, 1]$ to corresponding RGB colors.

EM2: for vector features like JJo and JJv, RGB channels are encoded based on XYZ values respectively as follows:

$$\overline{RGB}(h, w) = (\bar{f}_h^w - \min \bar{F}_h) / (\max(\bar{F}_h) - \min \bar{F}_h) \quad (7)$$

where $f_h^w \in R^3$ is the vector of h_{th} feature at w_{th} frame. Note that the operations are applied on each dimension.

EM3: this method encodes RGB channels based on scalar features from both subjects. Specifically, red channel is encoded based on features of main subject, green channel is encoded using features of the auxiliary subject, and blue channel is encoded based on both features. The encoding method is formulated as follows:

$$\begin{aligned} R(h, w) &= 1 - (f_h^w - \min F_h) / (\max(F_h) - \min F_h) \\ G(h, w) &= (v_h^w - \min V_h) / (\max(V_h) - \min V_h) \\ B(h, w) &= 4 \times R \times G \end{aligned} \quad (8)$$

where f, F and v, V represent features from main subject and other subject specifically.

2.4. CNN training and score fusion

In this paper, the CaffeNet (a version of Alexnet [13]) is adopted as the CNN model. The protocols used in [11, 12] are adopted to train the CNN models from scratch. Given a testing skeleton sequence, thirteen types of images are generated and each type of image is recognized with a trained CNN model. All the outputs (scores) of the CNN models are then fused into a final score by element-wise multiplication, which has been verified in [4, 11]. The fusion is done as follows:

$$\text{label} = F_{\max}(v_1 \circ v_2 \cdots v_{12} \circ v_{13}) \quad (9)$$

where v are the score vectors, \circ is the element-wise multiplication, and $F_{\max}(\cdot)$ is a function to find the index of the maximum element.

2.5. Implementation details

Joint coordinates are normalized in a way similar to the method in [12], where the spine lengths of the same subject in each frame (from 'base of the spine' to 'spine') are normalized to 1, and the other limb lengths are scaled in equal proportions. The scheme to select the main subject is adopted from [12], where the skeleton sequence having larger variations is set to be the main subject. Before selection, joint coordinates are translated from camera coordinate system to the body coordinate system, as described in [5]. The spatial features are directly calculated from normalized skeleton data to reduce the deviation introduced by the coordinate transformation.

Caffe was adopted as the CNN platform and a Nvidia Titan X GPU was used to run the experiments. The CNNs were trained using stochastic gradient descent (SGD) for a total of 30000 iterations. The models were trained from scratch and the weights were initialized using Gaussian Filter. The multi-step scheme was used to train the CNNs with step sizes as 10000, 18000, 24000, 28000 specifically. The learning rate was initially set to 0.01 and multiplied by 0.1 every epoch.

3. EXPERIMENT RESULTS

The proposed method was evaluated on NTU RGB+D Dataset. Currently, NTU RGB+D Dataset [12] is the largest dataset for action recognition. It has 56578 samples of 60 different actions classes, which are captured under 18 settings with different camera viewpoints and heights. The actions include single-subject cases and multi-subject interaction cases and are performed by 40 subjects aged between 10 and 35. This dataset is challenging and there are two types of protocols for evaluation of methods, cross-subject and cross-view. In this paper, the cross-view protocol is used. The effectiveness of different types of spatial features, different joint selection schemes were evaluated.

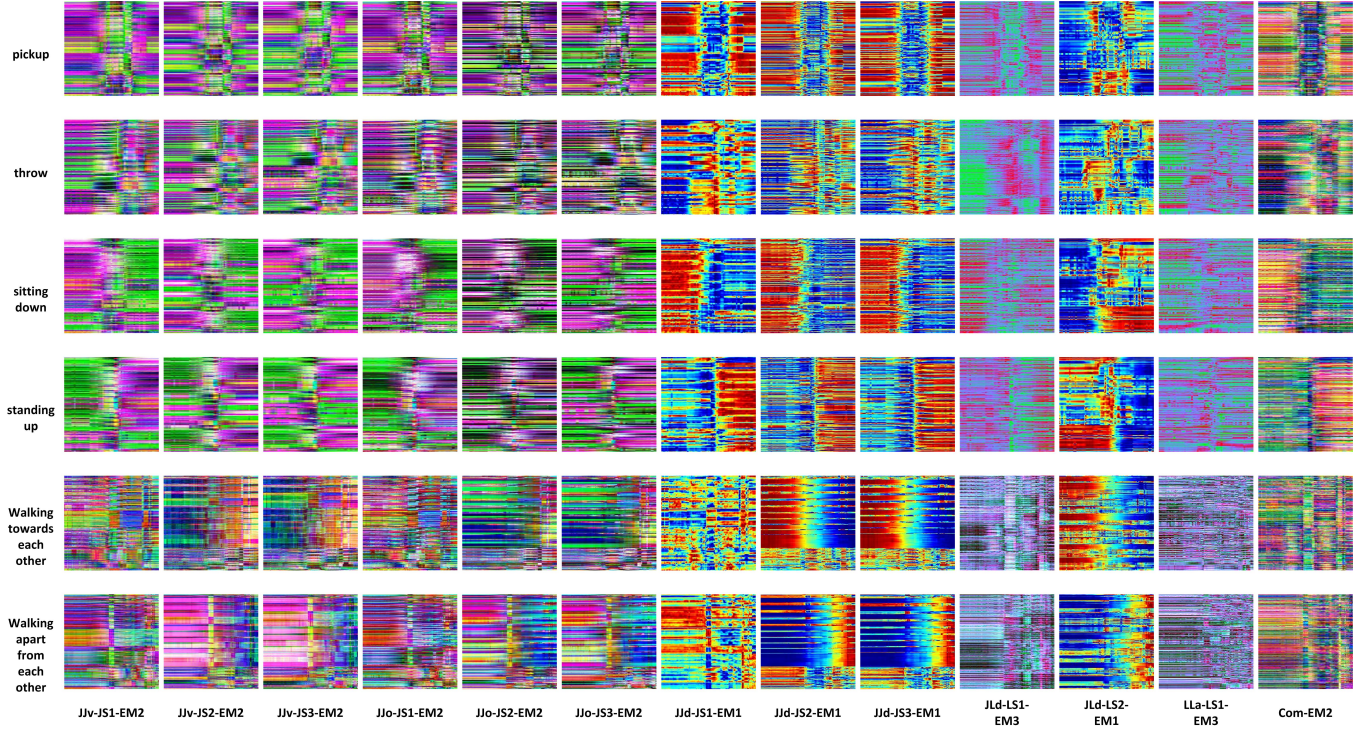


Fig. 2. Samples generated by the proposed method on NTU RGB+D Dataset. Six samples from different actions are visualized. The images in each row are generated from the same sample, and the images in each column are generated using the same method. The images within the same row represent the difference of methods, and the images within the same column represent the difference between action classes.

Table 1. Evaluation results of different features and encoding methods.

Feature	Method	Accuracy	Fused Accuracy	
JJv	JS1-EM2	62.45%	75.23%	82.31%
	JS2-EM2	65.12%		
	JS3-EM2	69.02%		
JJo	JS1-EM2	64.11%	73.51%	
	JS2-EM2 ²	55.14%		
	JS3-EM2	63.30%		
JJd	JS1-EM1	59.18%	73.01%	
	JS2-EM1	62.86%		
	JS3-EM1	62.95%		
JLd	LS1-EM3	63.08%	76.20%	
	LS2-EM1	59.71%		
LLa	LS1-EM3	62.57%	62.57%	
Com ¹	JS1&LS1-EM2	62.00%	62.00%	

Note 1: this method encodes the RGB channels based on JJd, JLd and LLa respectively. Note 2: this method is not used for final score fusion.

3.1. Evaluation of spatial features

The results of individual features and different encoding methods are listed in Table 1, as well as results of score-

Table 2. Experimental results (accuracy) on NTU RGB+D Dataset

Method	Accuracy
Lie Group[8]	52.76%
Dynamic Skeletons[14]	65.22%
HBRNN[15]	63.97%
Deep RNN[12]	64.09%
Part-aware LSTM[12]	70.27%
ST-LSTM+Trust Gate[16]	77.70%
JTM[9]	75.20%
Geometric Features[5]	82.39%
STA-LSTM[17]	81.20%
Proposed Method	82.31%

multiplication fusion. There are five features evaluated, each of which was evaluated with different feature (joint) selection methods and different encoding methods. The methods are denoted in the form ‘feature selection method - encoding method’, which have been described in Section 2.

As illustrated in Fig. 2, images generated from samples of different actions have discriminative textures. In addition, the spatial features are encoded into different textures by different methods.

From Table 1, it can be seen that the JJv feature is the best joint-joint feature, based on the comparisons of single results and fused results. Moreover, JLD seems to be the best feature among the five types of features, which coincides with the observations reported by [5]. Among the three kinds of joint selection methods, JS3 generally works better than the other two. This observation suggests that some of the joints are noise with regard to this task, which is consistent with the above analysis.

From Table 2 the results indicate that, compared with methods based hand-crafted features and those based on deep learning (RNNs and CNNs), the proposed method achieved state-of-the-art results.

4. CONCLUSIONS

In this paper, a method for skeleton-based action recognition using CNNs is proposed. This method explored encoding different spatial features into texture color images and achieved state-of-the-art results on NTU RGB+D Dataset. The experimental results indicated the effectiveness of texture images when used as spatio-temporal information representation, and the effectiveness of joint selection strategies for robust and cost-efficient computation.

5. REFERENCES

- [1] Wanqing Li, Zhengyou Zhang, and Zicheng Liu, "Action recognition based on a bag of 3D points," in *Proc. IEEE Conference on CVPR Workshops (CVPRW)*, 2010, pp. 9–14.
- [2] Pichao Wang, Wanqing Li, Zhimin Gao, Chang Tang, Jing Zhang, and Philip Ogunbona, "Convnets-based action recognition from depth maps through virtual cameras and pseudocoloring," in *Proc. ACM Conference on Multimedia*, 2015, pp. 1119–1122.
- [3] Pichao Wang, Wanqing Li, Zhimin Gao, Jing Zhang, Chang Tang, and Philip Ogunbona, "Action recognition from depth maps using deep convolutional neural networks," *IEEE Transactions on Human-Machine Systems*, vol. 46, no. 4, pp. 498–509, 2016.
- [4] Pichao Wang, Wanqing Li, Zhimin Gao, Yuyao Zhang, Chang Tang, and Philip Ogunbona, "Scene flow to action map: A new representation for RGB-D based action recognition with convolutional neural networks," in *CVPR*, 2017.
- [5] Songyang Zhang, Xiaoming Liu, and Jun Xiao, "On geometric features for skeleton-based action recognition using multilayer LSTM networks," in *Proc. IEEE Winter Conference on Applications of Computer Vision (WACV)*, 2017.
- [6] Lu Xia, Chia-Chih Chen, and JK Aggarwal, "View invariant human action recognition using histograms of 3D joints," in *Proc. IEEE Conference on CVPR Workshops (CVPRW)*, 2012, pp. 20–27.
- [7] Pichao Wang, Wanqing Li, Philip Ogunbona, Zhimin Gao, and Hanling Zhang, "Mining mid-level features for action recognition based on effective skeleton representation," in *Proc. IEEE Conference on Digital Image Computing: Techniques and Applications (DICTA)*, 2014, pp. 1–8.
- [8] Raviteja Vemulapalli, Felipe Arrate, and Rama Chellappa, "Human action recognition by representing 3d skeletons as points in a lie group," in *CVPR*, 2014, pp. 588–595.
- [9] Pichao Wang, Zhaoyang Li, Yonghong Hou, and Wanqing Li, "Action recognition based on joint trajectory maps using convolutional neural networks," in *Proc. ACM Conference on Multimedia*, 2016, pp. 102–106.
- [10] Yonghong Hou, Zhaoyang Li, Pichao Wang, and Wanqing Li, "Skeleton optical spectra based action recognition using convolutional neural networks," *IEEE Transactions on Circuits and Systems for Video Technology*, 2016.
- [11] Chuankun Li, Yonghong Hou, Pichao Wang, and Wanqing Li, "Joint distance maps based action recognition with convolutional neural network," *IEEE Signal Processing Letters*, 2017.
- [12] Amir Shahroudy, Jun Liu, Tian-Tsong Ng, and Gang Wang, "NTU RGB+D: A large scale dataset for 3D human activity analysis," in *CVPR*, June 2016.
- [13] Alex Krizhevsky, Ilya Sutskever, and Geoffrey E Hinton, "Imagenet classification with deep convolutional neural networks," *Advances in Neural Information Processing Systems*, pp. 1097–1105, 2012.
- [14] Eshed Ohn-Bar and Mohan Trivedi, "Joint angles similarities and hog2 for action recognition," in *Proc. IEEE Conference on CVPR Workshops (CVPRW)*, 2013, pp. 465–470.
- [15] Yong Du, Wei Wang, and Liang Wang, "Hierarchical recurrent neural network for skeleton based action recognition," in *CVPR*, 2015, pp. 1110–1118.
- [16] Jun Liu, Amir Shahroudy, Dong Xu, and Gang Wang, "Spatio-temporal lstm with trust gates for 3d human action recognition," in *ECCV*, 2016, pp. 816–833.
- [17] Sijie Song, Cuiling Lan, Junliang Xing, Wenjun Zeng, and Jiaying Liu, "An end-to-end spatio-temporal attention model for human action recognition from skeleton

data,” in *Proc. AAAI Conference on Artificial Intelligence*, 2017, pp. 4263–4270.