

SKELETON-BASED ACTION RECOGNITION WITH CONVOLUTIONAL NEURAL NETWORKS

Chao Li, Qiaoyong Zhong, Di Xie, Shiliang Pu

Hikvision Research Institute, Hangzhou, China
{lichao15, zhongqiaoyong, xiedi, pushiliang}@hikvision.com

ABSTRACT

Current state-of-the-art approaches to skeleton-based action recognition are mostly based on recurrent neural networks (RNN). In this paper, we propose a novel convolutional neural networks (CNN) based framework for both action classification and detection. Raw skeleton coordinates as well as skeleton motion are fed directly into CNN for label prediction. A novel skeleton transformer module is designed to rearrange and select important skeleton joints automatically. With a simple 7-layer network, we obtain 89.3% accuracy on validation set of the NTU RGB+D dataset. For action detection in untrimmed videos, we develop a window proposal network to extract temporal segment proposals, which are further classified within the same network. On the recent PKU-MMD dataset, we achieve 93.7% mAP, surpassing the baseline by a large margin.

Index Terms— Skeleton, CNN, Window Proposal Network, Action Recognition, Action Detection

1. INTRODUCTION

Articulated human pose, also referred to as skeleton, captures full information needed to understand the underlying activity of the subject. Compared with other modalities (e.g. RGB images, depth maps), skeleton data are more robust to noise like background and irrelevant objects. With the development of low-cost human skeleton capture systems (e.g. Kinect), large-scale 3D skeleton datasets have been made available [1, 2], which attract many research efforts [3] on skeleton-based human action recognition and detection. An ever-increasing use of skeleton data in a wide range of applications from human-computer interaction, virtual reality to video surveillance can be expected.

Considering the time series property of skeleton sequences in videos, recurrent neural networks (RNN), in particular long-short term memory networks (LSTM) are natural choices. Indeed the current state-of-the-art approaches are mostly based on LSTM. In this paper, we propose a new representation of skeleton data with convolutional neural networks (CNN), which is shown to outperform a strong LSTM baseline. Besides, we adapt the widely used Faster

R-CNN [4] object detection framework to action detection in temporal domain. With the novel CNN-based detection framework, we obtain 58% *absolute* mAP improvement over the baseline.

2. RELATED WORKS

2.1. Representation of skeleton

LSTM has been well exploited to model the temporal pattern of skeleton sequences. Within the LSTM framework, many improvements have been made in the literature. For example, [5] explored the co-occurrence feature of skeleton joints. [6] exploited attention model in both spatial and temporal domain. Recently, [7] developed a view adaptive RNN to cope with the viewpoint variations explicitly. On the other hand, Ke et al. [8] proposed a CNN based representation of skeleton and achieved state-of-the-art performance. Our CNN representation differs from [8] on the input form of skeleton as well as the network architecture. Besides, we get significantly superior performance over [8].

2.2. CNN for object detection

CNN has achieved great success in many image recognition tasks, e.g. image recognition, object detection. For object detection, Faster R-CNN [4] is the current state of the art. It consists of two cascaded stages. In stage 1, a fully convolutional region proposal network (RPN) is utilized to extract putative region proposals. In stage 2, features of the proposals are ROI-pooled and further classified with R-CNN. By sharing features between RPN and R-CNN, real-time detection is achieved.

Faster R-CNN was originally designed for object detection in still images. [9] adapted the framework to temporal activity detection in RGB videos with a 3D convolutional network. In this work, we are the first to adapt Faster R-CNN to the task of skeleton-based temporal action detection.

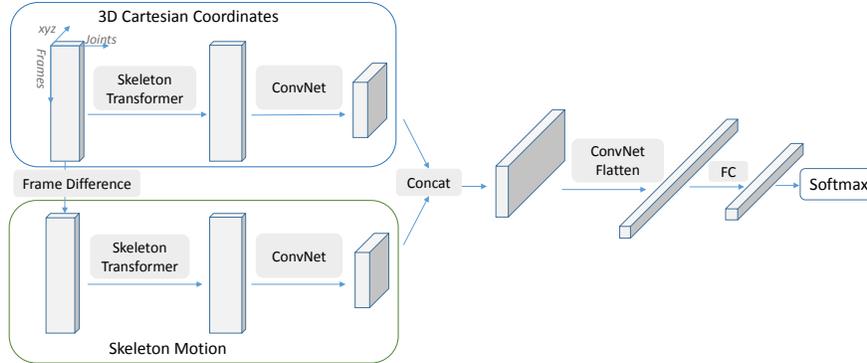


Fig. 1. CNN representation of skeleton sequences for action classification.

3. METHOD

3.1. Action classification

We design a simple yet effective convolutional architecture (Fig. 1) for action classification from trimmed skeleton sequences. Besides raw joint coordinates, motion of skeleton joints from two consecutive frames are fed as an extra input to the network. In addition, we propose a novel skeleton transformer module. With the module, the network is able to automatically learn a better ordering of joints as well as new joints that are more informative than arbitrarily given ones. To deal with multi-person settings, we use maxout to merge features from skeletons of different individuals.

3.1.1. Two-stream CNN

Two-stream architecture was first introduced in [10], where RGB images and optical flow fields are utilized as two input streams of a network. Similarly, we define two network inputs for the case of skeleton data. Given a 3D joint coordinate $\mathbf{J} = (x, y, z)$, skeleton of one person is represented as a set of joint coordinates $\mathbf{S} = \{\mathbf{J}_1, \mathbf{J}_2, \dots, \mathbf{J}_N\}$ where N is the number of joints per skeleton. Skeleton motion between two consecutive frames is computed as $\mathbf{M} = \mathbf{S}^{t+1} - \mathbf{S}^t = \{\mathbf{J}_1^{t+1} - \mathbf{J}_1^t, \mathbf{J}_2^{t+1} - \mathbf{J}_2^t, \dots, \mathbf{J}_N^{t+1} - \mathbf{J}_N^t\}$ where t is frame index. A skeleton sequence of T frames can be represented as a $T \times N \times 3$ array, which is treated as a $T \times N$ sized 3-channel image. Raw skeleton coordinates \mathbf{S} and skeleton motion \mathbf{M} are used as two input streams of our network. Note that we do not perform any normalization.

Actions in a video may span varying length of frames. Since skeleton data are treated as image data, we normalize all videos to a fixed length by a simple image resizing operation.

3.1.2. Skeleton transformer

For an image, the semantic continuity of pixels is critical. For example, if we shuffle the location of all pixels randomly, the resulting image would be non-sense and difficult to recognize

for both humans and machines. For the $T \times N \times 3$ skeleton image data, ordering of N joints are arbitrarily chosen (e.g. left eye, right eye, nose, ...), which may not be optimal. To address this issue, we propose a skeleton transformer module. Given an $N \times 3$ skeleton \mathbf{S} , we perform a linear transformation $\mathbf{S}' = (\mathbf{S}^T \cdot \mathbf{W})^T$, where \mathbf{W} is an $N \times M$ weight matrix. \mathbf{S}' is a list of M new interpolated joints. Note that both ordering and location of the joints are rearranged. The network selects important body joints automatically, which can be interpreted as a simple variant of attention mechanism.

Skeleton transformer can be implemented simply with a fully connected layer (without bias). We place this module at the very beginning of the network before convolution layers such that it is trained end to end.

3.1.3. Multi-person maxout

The methods mentioned above are designed for the case of single person. For those activities involving human-human interaction (e.g. hugging, shaking hands), there will be multiple people. A common choice in the literature is to concatenate skeletons of different people as the network input. Zero padding is required to deal with varying number of people.

In this paper, we adopt the maxout [11] scheme for multiple people. Skeletons of different people go through the same network layers, and their feature maps are merged by an element-wise maximum operation after the last convolution layer. The advantage is two-fold. Firstly, the varying number of people issue can be resolved gracefully without zero padding. Secondly, by weight sharing, our method can be extended from two people to more people without increasing model size.

3.1.4. Network architecture

We design a tiny 7-layer network which consists of 3 convolution layers and 4 fully connected layers (at which point performance saturates). Our network contains only 1.3 million parameters. And it can be easily trained from scratch with-

out any pre-training. Compared with [8], where an ImageNet pre-trained VGG19 net is used, our model is superior on its compact model size and fast inference speed as well.

3.2. Action detection

By interpreting a sequence of skeleton data as a $T \times N \times 3$ image, it is straightforward to adapt object detection methods to the task of action detection in temporal domain. In this paper, we take Faster R-CNN as an example. Other object detection frameworks should work as well.

As displayed in Fig. 2, the region proposal network (RPN) is replaced with a window proposal network (WPN). In particular, 2D anchors are flattened to 1D anchors. Window proposals along the temporal dimension are extracted based on pre-defined anchors. Window regression instead of bounding box regression is performed to refine the temporal position of window proposals. After the proposals are ready, we pool features of each window from the shared feature maps with the crop-and-resize operation. These features are then fed to the R-CNN subnetwork for classification and window regression. For the backbone network, we use the same architecture as action classification in Fig. 1.

In our experiments, we use 4 anchor scales, i.e. {50, 100, 200, 400}. For other hyper-parameters, we follow the settings recommended in [4]. During training, we randomly choose a temporal scale factor between 0.8 and 1.5. During testing, we use single scale (the original resolution).

4. EXPERIMENTS

We validate our method on two large-scale skeleton datasets. The NTU RGB+D dataset [1] is designed for action classification task. It contains 56880 well trimmed video clips spanning 60 action categories. The very recent PKU-MMD dataset [2] is designed for action detection task, which contains 1076 untrimmed videos and 21545 action instances. The number of action categories is 51.

4.1. Action classification

4.1.1. Ablation study

To evaluate contributions of different components, we perform an ablation study on the NTU RGB+D dataset. Table 1 shows the results. Using plain CNN, we already outperform STA-LSTM [6] (see Table 2), a strong LSTM baseline. Skeleton motion improves accuracy by 1.6 and 3.3 points in cross-subject and cross-view settings respectively. Skeleton transformer improves cross-subject by 1.8 points, while improvement on cross-view is marginal. This could be explained by the reason that variation of actions across subjects is larger than that of across views, which can be alleviated by skeleton

Table 1. Ablation study on action classification. Accuracies are measured on validation set of the NTU RGB+D dataset.

| Method | Cross-subject | Cross-view |
|------------------|---------------|--------------|
| CNN | 0.798 | 0.852 |
| CNN+Motion | 0.814 | 0.885 |
| CNN+Trans | 0.816 | 0.854 |
| CNN+Motion+Trans | 0.832 | 0.893 |

Table 2. Action classification performance on validation set of the NTU RGB+D dataset.

| Method | Cross-subject | Cross-view |
|---------------|---------------|--------------|
| STA-LSTM [6] | 0.734 | 0.812 |
| VA-LSTM [7] | 0.792 | 0.877 |
| Ke et al. [8] | 0.796 | 0.848 |
| Proposed | 0.832 | 0.893 |

transformer. Combining skeleton motion and skeleton transformer, we obtain 83.2% and 89.3% accuracy in the two partitioning schemes respectively.

4.1.2. Comparison to the state-of-the-arts

Our method significantly outperforms all recent state-of-the-art approaches in both cross-subject and cross-view settings (Table 2). Specifically, we improve accuracy by 10 points over STA-LSTM in cross-subject setting. Besides, our method is also superior over Ke et al. [8], which is also based on CNN. The excellent result clearly proves the ability of CNN to model temporal pattern. We believe that CNN can be applied to other time series signals other than skeleton sequences.

4.2. Action detection

We validate our action detection pipeline on the PKU-MMD dataset. Since the official evaluation code is unavailable, we report our performance based on our own implementation of mean average precision (mAP) over different actions. Table 3 shows the mAP numbers at two different Intersection over Union (IoU) thresholds. The detector easily enjoys benefits from the capability of our CNN-based classifier. Compared with the strong baseline JCRRNN [12], we obtain a performance boost. That is 58% *absolute* mAP improvement in cross-subject and 40% mAP improvement in cross-view at IoU threshold of 0.5. The performance improvement indicates that it is a viable solution to treat skeleton sequences as images and transform the temporal action detection problem into a unidimensional object detection problem.

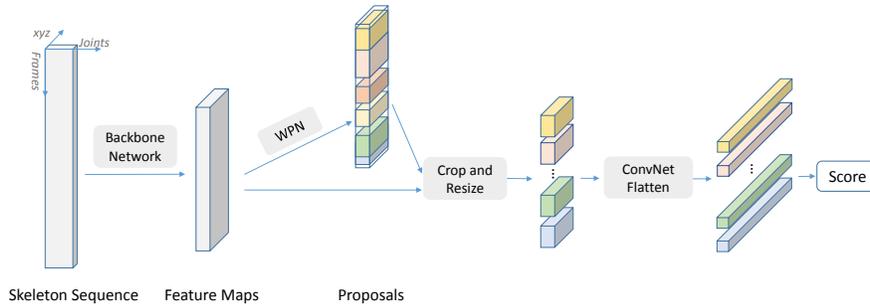


Fig. 2. Skeleton-based temporal action detection pipeline.

Table 3. Action detection performance on validation set of the PKU-MMD dataset. θ is the IoU threshold.

| Partition | Cross-subject | | Cross-view | |
|-------------|---------------|--------------|--------------|--------------|
| | θ | θ | θ | θ |
| JCRRNN [12] | 0.452 | 0.325 | 0.699 | 0.533 |
| Proposed | 0.922 | 0.904 | 0.958 | 0.937 |

5. CONCLUSION

Skeleton based human action recognition is drawing more and more attention as the popularity of 3D skeleton data. By treating skeleton sequences as images, we propose a novel CNN based framework for both action classification and detection tasks. Our method achieves new state-of-the-art performance on two recent large-scale skeleton datasets. The proposed action detection approach detects actions in a batch-processing way, while online detection is required for real-time applications. We leave it as future exploration.

6. REFERENCES

- [1] Amir Shahroudy, Jun Liu, Tian-Tsong Ng, and Gang Wang, “NTU RGB+D: A large scale dataset for 3d human activity analysis,” in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2016, pp. 1010–1019.
- [2] Chunhui Liu, Yueyu Hu, Yanghao Li, Sijie Song, and Jiaying Liu, “PKU-MMD: A large scale benchmark for continuous multi-modal human action understanding,” *arXiv preprint arXiv:1703.07475*, 2017.
- [3] Fei Han, Brian Reily, William Hoff, and Hao Zhang, “Space-time representation of people based on 3d skeletal data: A review,” *Computer Vision and Image Understanding*, 2017.
- [4] Shaoqing Ren, Kaiming He, Ross Girshick, and Jian Sun, “Faster R-CNN: Towards real-time object detection with region proposal networks,” in *Advances in neural information processing systems*, 2015, pp. 91–99.
- [5] Wentao Zhu, Cuiling Lan, Junliang Xing, Wenjun Zeng, Yanghao Li, Li Shen, and Xiaohui Xie, “Co-occurrence feature learning for skeleton based action recognition using regularized deep lstm networks,” *arXiv preprint arXiv:1603.07772*, 2016.
- [6] 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 *AAAI Conference on Artificial Intelligence*, 2017.
- [7] Pengfei Zhang, Cuiling Lan, Junliang Xing, Wenjun Zeng, Jianru Xue, and Nanning Zheng, “View adaptive recurrent neural networks for high performance human action recognition from skeleton data,” *arXiv preprint arXiv:1703.08274*, 2017.
- [8] Qihong Ke, Mohammed Bennamoun, Senjian An, Ferdous Sohel, and Farid Boussaid, “A new representation of skeleton sequences for 3d action recognition,” *arXiv preprint arXiv:1703.03492*, 2017.
- [9] Huijuan Xu, Abir Das, and Kate Saenko, “R-C3D: Region convolutional 3d network for temporal activity detection,” *arXiv preprint arXiv:1703.07814*, 2017.
- [10] Karen Simonyan and Andrew Zisserman, “Two-stream convolutional networks for action recognition in videos,” in *Advances in neural information processing systems*, 2014, pp. 568–576.
- [11] Ian J Goodfellow, David Warde-Farley, Mehdi Mirza, Aaron C Courville, and Yoshua Bengio, “Maxout networks,” *ICML*, vol. 28, pp. 1319–1327, 2013.
- [12] Yanghao Li, Cuiling Lan, Junliang Xing, Wenjun Zeng, Chunfeng Yuan, and Jiaying Liu, “Online human action detection using joint classification-regression recurrent neural networks,” in *European Conference on Computer Vision*. Springer, 2016, pp. 203–220.