

# Feedback Graph Convolutional Network for Skeleton-based Action Recognition

Hao Yang<sup>1</sup>, Dan Yan<sup>1</sup>, Li Zhang<sup>1,2</sup>, Dong Li<sup>1,2</sup>, YunDa Sun<sup>1</sup>, Shaodi You<sup>3</sup>,  
and Stephen J. Maybank<sup>4</sup>

<sup>1</sup> R&D Center of Artificial Intelligent, NUCTECH Company Limited, Beijing, China  
{yanghao1, yandan, sunyunda, li.dong}@nuctech.com

<sup>2</sup> Department of Engineering Physics, Tsinghua University, Beijing, China  
zli@mail.tsinghua.edu.cn

<sup>3</sup> Informatics Institute, University of Amsterdam. Amsterdam, Netherlands  
s.you@uva.nl

<sup>4</sup> Department of Computer Science and Information Systems, Birkbeck College,  
London, United Kingdom  
sjmaybank@dcs.bbk.ac.uk

**Abstract.** Skeleton-based action recognition has attracted considerable attention in computer vision since skeleton data is more robust to the dynamic circumstance and complicated background than other modalities. Recently, many researchers have used the Graph Convolutional Network (GCN) to model spatial-temporal features of skeleton sequences by an end-to-end optimization. However, conventional GCNs are feedforward networks which are impossible for low-level layers to access semantic information in the high-level layers. In this paper, we propose a novel network, named Feedback Graph Convolutional Network (FGCN). This is the first work that introduces the feedback mechanism into GCNs and action recognition. Compared with conventional GCNs, FGCN has the following advantages: (1) a multi-stage temporal sampling strategy is designed to extract spatial-temporal features for action recognition in a coarse-to-fine progressive process; (2) A dense connections based Feedback Graph Convolutional Block (FGCB) is proposed to introduce feedback connections into the GCNs. It transmits the high-level semantic features to the low-level layers and flows temporal information stage by stage to progressively model global spatial-temporal features for action recognition; (3) The FGCN model provides early predictions. In the early stages, the model receives partial information about actions. Naturally, its predictions are relatively coarse. The coarse predictions are treated as the prior to guide the feature learning of later stages for a accurate prediction. Extensive experiments on the datasets, NTU-RGB+D, NTU-RGB+D120 and Northwestern-UCLA, demonstrate that the proposed FGCN is effective for action recognition. It achieves the state-of-the-art performance on the three datasets.

**Keywords:** Feedback, Graph Convolutional Network, Skeleton, Action Recognition

## 1 Introduction

In recent years, the quantity of videos uploaded from various terminals has exploded. This has driven imperious demands for human action analysis automatically based on the content of videos. In particular, human action recognition using skeleton has attracted many computer vision researchers because of its strong adaptability to the effects of dynamic circumstance and complicated background, as compared with other modalities such as RGB [16] and optical flow [47]. Early deep learning methods using skeletons for action recognition usually represent the skeleton data as a sequence of joint-coordinate vectors [6,48,58,24] or a pseudo-image [17,32,27] which is then modeled by a RNN or CNN respectively. However, these methods do not explicitly exploit the spatial dependencies among correlated joints, even though the spatial dependencies are informative for understanding human actions. More recently, some methods [56,44,43,23] construct spatial temporal graphs based on the natural connections of joints and temporal edges of consecutive frames. They then exploit a GCN to model spatial-temporal features. However, the conventional GCNs are all single-pass feedforward networks that are fed with the entire skeleton sequence. It is difficult for these methods to extract effective spatial-temporal features, because the useful information is usually buried in the motion-irrelevant or undiscriminating clips when they are fed with entire skeleton sequence. For example, in the action “kicking something”, most clips are “standing upright”, and in the action “wear a shoe”, most clips are a subject sitting on a chair. Then the single-pass feedforward networks can not access the high-level semantic information for the low-level layers. Meanwhile, inputting the entire skeleton sequence increases computational complexity of the model.

Motivated by this, we propose a novel neural network, named Feedback Graph Convolutional Network (FGCN), to extract effective spatial-temporal features from skeleton data in a coarse-to-fine progressive process for action recognition. The FGCN is the first work that introduces feedback mechanism into GCNs and action recognition. Compared with conventional GCNs, the FGCN has a multi-stage temporal sampling strategy which divides input skeleton sequences into multiple stages in the temporal domain and sparsely samples input skeleton clips from temporal stages to avoid feeding with the entire skeleton sequence. Each sampled clip is input into graph convolutional layers to extract local spatial-temporal features for each stage. A Feedback Graph Convolutional Block (FGCB) is proposed to model global spatial-temporal features by fusing the local features. The FGCB is a local dense graph convolutional network with lateral connections from each stage to the next stage and it introduces feedback connections into conventional GCNs. From a semantic point of view it works in a top down manner, which makes it possible for low-level convolutional layers to access semantic information in the high-level layers at each stage. From the temporal domain, the feedback mechanism in FGCB works with a sequence of cause-and-effect and the output of the previous stage flows into the next stage to modulate its input.

Another advantage of the FGCN is that it provides early predictions of the output in a fraction of the total inference time. This is valuable in many applications such as robotics or autonomous driving, in which latency time is very crucial. The early predictions are a result of the proposed multi-stage coarse-to-fine progressively optimization. In the early stages, FGCN is only fed with a part of skeleton sequence and the information about the action is limited, so the inferences of it are relatively coarse. These inferences are treated as a prior to guide the feature learning in later stages. In later stages, the model receives more complete information about the action and the guider of former inferences, thus it outputs more accurate inferences. Several temporal fusion strategies are proposed to fuse the local predictions in temporal stages for a video-level prediction. The strategies enable the network to be optimized in a progressive process.

The main contributions of this paper are summarized as follows:

- We propose a novel Feedback Graph Convolutional Network (FGCN) for action recognition from skeleton sequences. It models spatial-temporal features by a multi-stage progressive process. To our knowledge, this is the first work that introduces the feedback mechanism into GCNs and action recognition.
- We propose a dense connections based Feedback Graph Convolutional Block (FGCB) which is a local network with lateral connections between two temporal stages. Functionally, it transmits high-level semantic features as priors to module its features in low-level layers.
- The FGCN model provides early predictions, which benefits from the multi-stage coarse-to-fine progressive optimization. The proposed model is extensively evaluated on three datasets, NTU-RGB+D, NTU-RGB+D120 and Northwestern-UCLA, and it achieves state-of-the-art performance on the three datasets.

## 2 Related Works

### 2.1 Skeleton based Action Recognition

As the depth sensor technologies (*i.e.* kinect [59]) and pose estimation algorithms [52,4] matured, it becomes possible to capture skeleton data in real time by locating the key joints. The skeleton data is robust to illumination change, scene variation, and complex background. These facilitate the data-driven method’s development of skeleton-based action recognition. Conventional action recognition methods usually extract hand-crafted features from skeleton sequences. Some traditional methods [8,53,33,40] design several view-invariant features of actions. Examples of these features are body part-based skeletal quads [8,53], group sparsity based class-specific dictionary coding [33], and canonical view transformed features [40]. Other traditional methods integrate the information from different modalities that are always available in 3D action datasets. Some works [13,35,39,54] combine the depth information with the skeleton to improve performance. The depth information is represented by HOG features [13,35] and Fourier Temporal Pyramids [54], or it is modeled by random decision forests

[39]. The recent success of deep learning has led to a surge of deep network based skeleton modeling methods. The widely used models are RNNs and CNNs. RNN-based methods [6,48,58,24] usually concatenate all of the joint-coordinates (2D or 3D) in each frame as a vector and then model the features of actions by a RNN fed with a sequence of the coordinate vectors. CNN-based methods [17,32,27] stack the sequence of coordinate vectors to obtain a pseudo-image, and then reduce the action recognition using skeleton sequences to an image classification task. The two-stream based model [60] combines RNN and CNN, operating on coordinate vectors of skeletons and RGB images respectively, to improve performance from a single network. However, these methods do not explicitly model the spatial dependence between correlated joints which is crucial for understanding human actions.

## 2.2 GCN based Action Recognition

The Graph Convolutional Networks (GCNs) [2,34,7,12,19] generalize the convolutional operation to deal with the data with graph construction. There are two main ways of constructing GCNs: spatial perspective and spectral perspective. Spatial perspective methods [2,34] directly perform the convolution filters on the graph vertexes and their neighbors. In contrast, spectral perspective methods [7,12,19] consider the graph convolution as a form of spectral analysis by utilizing the eigenvalues and eigenvectors of the graph Laplacian matrices. This work follows the spatial perspective based methods [56,44,43,23]. The ST-GCN model [56] is proposed to move beyond the limitations of hand-crafted parts and traversal rules used in previous methods. It operates on a spatial temporal graph to model the structured information about the joints along both the spatial and temporal dimensions. Based on ST-GCN, the 2s-AGCN model [44] proposes a two-stream adaptive graph convolutional network, which exploits the second-order information of the skeleton to improve the performance of action recognition. The DGNN model [43] represents the skeleton data as a directed acyclic graph based on the kinematic dependency between the joints and bones. The AS-GCN model [23] proposes an actional-structural graph convolution network by generating the skeleton graph with actional links and structural links. However, conventional GCNs are all feedforward networks in which it is impossible for low-level layers to access the semantic information in high-level layers.

## 2.3 Feedback Network

Feedback mechanism exists in the human visual cortex [15,9], and it has been a focus of research in psychology [1] and control theory [20,37]. In recent years, feedback mechanism has been introduced into deep neural networks in computer vision [49,57,26,11,10,5], because it allows the network to carry the information of output to correct previous states. In object recognition, the dasNet model [49] exploits the feedback structure by dynamically altering its convolutional filter sensitivities during classification and iteratively focusing its internal attention on some of its convolutional filters. Feedback Network [57] firstly introduces

the feedback mechanism into the convolutional recurrent neural network, which transfers the hidden state with high-level information to the input layer. In super resolution, several efforts [26,11,10] are made to take advantage of the feedback mechanism. The DBPN model [11] proposes a deep back-projection network which exploits iterative up-projection and down-projection units to achieve error feedback. The DSRN model [10] proposes a dual-state RNN and transmits the information between two recurrent states via a delayed feedback. The SRFBN model [26] designs a feedback block to handle the feedback connections and refines low-level representations with high-level information. In human pose estimation, [5] proposes an iterative error feedback (IEF) by iteratively estimating and applying a self-correction to the current estimation.

### 3 The Method

#### 3.1 Graph Convolutional Network

GCNs generalize the convolution operation to learn effective representations from graph structured data. In action recognition, the skeleton of a body is defined as an undirected graph in which each joint of the skeleton is defined as a vertex of the graph and the natural connections in the human body are defined as edges of the graph. In this paper, the skeleton in the frame  $t$  is denoted as a graph  $G_t = \{\mathbf{V}_t, \mathbf{E}_t\}$ , where  $\mathbf{V}_t$  is the set of joints in the frame and  $\mathbf{E}_t$  is the set of bones in the skeleton. For 3D skeleton data, the joint set is denoted as  $\mathbf{V}_t = \{v_{ti}\}_{i=1}^N$ , where  $v_{ti} = (x_{ti}, y_{ti}, z_{ti})$ . Given two joints  $v_{ti} = (x_{ti}, y_{ti}, z_{ti})$  and  $v_{tj} = (x_{tj}, y_{tj}, z_{tj})$ , a bone of the skeleton is defined as a vector  $e_{v_{ti}, v_{tj}} = (x_{tj} - x_{ti}, y_{tj} - y_{ti}, z_{tj} - z_{ti})$ ,  $(i, j) \in Q$ , where  $Q$  is the set of naturally connected human body joints. The skeleton sequence with  $len$  frames is denoted as  $S = \{G_1, G_2, \dots, G_{len}\}$ .

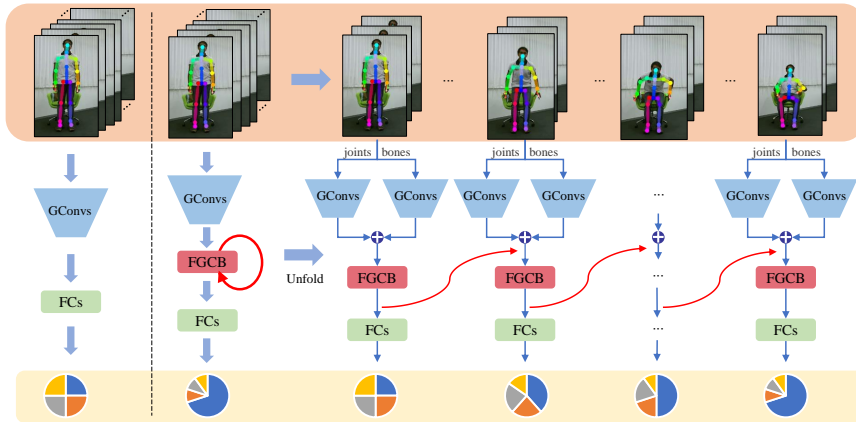
The graph convolution is defined operating on each vertex and its neighbors. For a vertex  $v_{ti}$  in the graph, its neighbor set is denoted as  $N(v_{ti}) = \{v_{tj} | d(v_{ti}, v_{tj}) \leq D\}$ , where  $d(v_{ti}, v_{tj})$  is the length of the shortest path from  $v_{tj}$  to  $v_{ti}$ . We set  $D = 1$  for the 1-distance neighbor set in this paper. The graph convolution operating on the neighbor set of vertex  $v_{ti}$  is formulated as:

$$f_{out}(v_{ti}) = \sum_{v_{tj} \in N(v_{ti})} \frac{1}{Z[l(v_{tj})]} f_{in}(v_{tj}) W[l(v_{tj})], \quad (1)$$

where  $f_{in}$  and  $f_{out}$  denote the input and output feature maps of this convolutional layer.  $l(v_{tj})$  is the label function which allocates a label from 1 to  $K$  for the vertex in  $N(v_{ti})$ . In our experiments, we set  $K = 3$  empirically to divide  $N(v_{ti})$  into 3 subsets.  $W(\cdot)$  is the weighting function which provides a weight vector according to the label  $l(v_{tj})$ . Similarly,  $Z[l(v_{tj})]$  denotes the number of vertices corresponding to the subset of  $l(v_{tj})$ .

In implementation, the connections of a graph are recorded in an  $N \times N$  adjacency matrix  $\mathbf{A}_k$ . With the adjacency matrix, Eqn. 1 can be formulated as:

$$f_{out} = \sum_{k=1}^K \mathbf{W}_k (\Lambda_k^{-\frac{1}{2}} \mathbf{A}_k \Lambda_k^{-\frac{1}{2}} f_{in}) \odot (\mathbf{M}_k), \quad (2)$$



**Fig. 1.** Comparison of the conventional GCNs (left) and the proposed FGCN (right). Red arrows represent the feedback connections of the feedback block (FGCB).

where  $\odot$  denotes the dot product and  $\Lambda_k^{ii} = \sum_j \mathbf{A}_k^{ij}$  is a diagonal matrix.  $\mathbf{W}_k$  is the weight vector of the convolution operation, which corresponds to the weighting function  $W(\cdot)$  in Eqn. 1. In practice,  $\mathbf{A}_k$  is allocated with a learnable weight matrix  $\mathbf{M}_k$  which is an  $N \times N$  attention map that indicates the importance of each vertex. It is initialized as an all-one matrix.

### 3.2 Feedback Graph Convolutional Network

Traditional action recognition methods [56,44,43,23] based on GCNs are all fed with the entire skeleton sequence in a feedforward network. However, the useful information is usually buried in the motion-irrelevant and undiscriminating clips when fed with entire skeleton sequence. And single-pass feedforward networks can not access semantic information at low-level layers. To tackle these problems, we propose a Feedback Graph Convolutional Network (FGCN) which extracts spatial-temporal features by a multi-stage progressive process, as shown in Fig. 1. Specifically, in the FGCN a multi-stage temporal sampling strategy is designed to sparsely sample a sequence of input clips from the skeleton data, instead of operating on the entire skeleton sequence directly. These clips are first fed into graph convolutional layers to extract the local spatial-temporal features. Then, a Feedback Graph Convolutional Block (FGCB) is proposed to fuse the local spatial-temporal features from multiple temporal stages by transmitting the high-level information in the previous stage to the next stage to modulate its input. Finally, several temporal fusion strategies are proposed to fuse the local predictions from all temporal stages to give a video-level prediction.

Formally, given a skeleton sequence  $S$ , the multi-stage temporal sampling strategy first divides it into  $T$  temporal stages with equal time interval, denoted as  $S = \{s_1, s_2, \dots, s_T\}$ . In each temporal stage, a skeleton clip is sampled randomly as an input of the deep model, denoted as  $\{c_1, c_2, \dots, c_T\}$ , where  $c_t$  is

the input clip sampled from the corresponding stage  $s_t$ . Each sampled clip  $c_t$  is input to the stacked multiple graph convolutional layers to extract the local spatial-temporal features in the corresponding temporal stage, formulated as:

$$\mathbf{F}_t = f_{GConv}(c_t), \quad (3)$$

where  $t = 1, 2, \dots, T$ , and  $\mathbf{F}_t$  is the local spatial-temporal features extracted by graph convolutional layers which are denoted as  $GConv$  in Fig. 1.

The local features extracted from all temporal stages flow into the feedback block FGCB to learn global spatial-temporal features for action recognition. As shown in Fig. 2, FGCB receives two inputs at the stage  $t$ : one is the hidden state from the previous stage  $t - 1$ , denoted as  $\mathbf{H}_{t-1}$ ; the other is the local features from the current stage, denoted as  $\mathbf{F}_t$ . Particularly, the input feature at the first stage  $\mathbf{F}_1$  is regarded as the initial hidden state  $\mathbf{H}_0$ . Based on these two inputs, the feedback process of FGCB is formulated as:

$$\mathbf{H}_t = f_{FGCB}(\mathbf{H}_{t-1}, \mathbf{F}_t), \quad (4)$$

where  $\mathbf{H}_t$  is the output of FGCB at stage  $t$ , and the function  $f_{FGCB}(\cdot)$  represents the operations of the feedback block FGCB. More details about FGCB can be found in Section 3.3.

Following the FGCB, a fully connected layer and a softmax loss layer are used at each stage to predict actions. The prediction process from the output  $H_t$  of FGCB is formulated as:

$$P_t = f_{pred}(\mathbf{H}_t), \quad (5)$$

where  $P_t \in R^C$  denotes the local prediction at stage  $t$  and  $C$  is the number of actions. The function  $f_{pred}(\cdot)$  represents the operations of the fully connected layer and the softmax layer. After operating on  $T$  temporal stages, we will obtain totally  $T$  local predictions  $\{P_1, P_2, \dots, P_T\}$ . Several temporal fusion strategies are proposed to fuse these local predictions corresponding to multiple stages for a video-level prediction  $P_S$  which is computed as:

$$P_S = f_{tf}(P_1, P_2, \dots, P_T), \quad (6)$$

where  $f_{tf}$  is the operations of a temporal fusion strategy. In this paper, we propose three temporal fusion strategies, *i.e.* last-win-all fusion, average fusion and weighting fusion. The FGCN model is trained end-to-end with the cross-entropy loss as follows:

$$L(y, P_S) = - \sum_{i=1}^C y^i \log(P_S^i), \quad (7)$$

where  $y$  is the action label of the skeleton  $S$ , if  $y = i$ ,  $y^i$  is set as 1, otherwise it is set as 0.

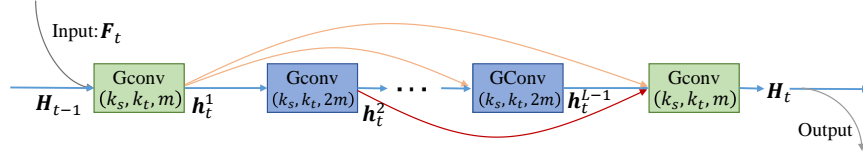


Fig. 2. The detailed architecture of the proposed FGCB local network.

### 3.3 Feedback Graph Convolutional Block

The feedback block FGCB is the core component of the FGCN model. On the one hand, the FGCB transmits the high-level semantic information back to low-level layers to refine their encoded features. On the other hand, the output at the previous stage flows into the next stage to modulate its input. To enable the FGCB to effectively transmit information from high-level to low-level and from the previous stage to the next stage, we propose a dense connected local graph convolutional network which adds shortcut connections from each layer to all subsequent layers. At a temporal stage  $t$ , the FGCB receives the high-level information from the output  $\mathbf{H}_{t-1}$  of the previous stage to modulate the low-level input  $\mathbf{F}_t$  of the current stage. In our model, the FGCB consists of  $L$  spatial temporal graph convolutional layers. The spatial temporal graph convolutional layer is denoted as  $GConv(k_s, k_t, m)$  in Fig. 2, where  $k_s$  and  $k_t$  are the kernel size in the spatial and temporal domains respectively, and  $m$  denotes output channels of the graph convolutional layer.

As shown in Fig. 2, the first convolutional layer in FGCB receives two inputs  $\mathbf{F}_t$  and  $\mathbf{H}_{t-1}$ . It compresses and fuses the features from the concatenation of the two inputs  $[\mathbf{F}_t, \mathbf{H}_{t-1}]$ . The output of this layer is formulated as:

$$\mathbf{h}_t^1 = f_{FGCB}^1([\mathbf{F}_t, \mathbf{H}_{t-1}]), \quad (8)$$

where  $f_{FGCB}^1(\cdot)$  denotes the operations in the first convolutional graph layer of FGCB, and  $\mathbf{h}_t^1$  denotes the output feature maps of the first layer. Following the first layer, the  $l_{th}$  layer receives the output feature maps from all preceding layers,  $\mathbf{h}_t^1, \mathbf{h}_t^2, \dots, \mathbf{h}_t^{l-1}$ , as input:

$$\mathbf{h}_t^l = f_{FGCB}^l([\mathbf{h}_t^1, \mathbf{h}_t^2, \dots, \mathbf{h}_t^{l-1}]), \quad (9)$$

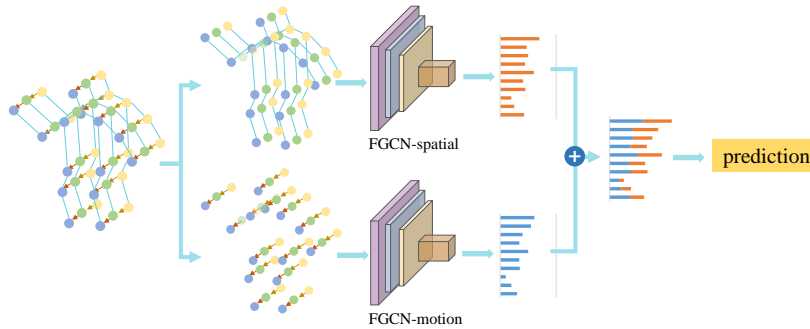
where  $l = 1, 2, \dots, L$  and  $[\mathbf{h}_t^1, \mathbf{h}_t^2, \dots, \mathbf{h}_t^{l-1}]$  refers to the concatenated feature maps in preceding layers. Similar to the first layer, the final layer in FGCB compresses and fuses the feature maps from the outputs of all preceding layers to produce the output of FGCB:

$$\mathbf{H}_t = \mathbf{h}_t^L = f_{FGCB}^L([\mathbf{h}_t^1, \mathbf{h}_t^2, \dots, \mathbf{h}_t^{L-1}]), \quad (10)$$

### 3.4 Two-stream Framework of FGCN

The joints and bones of a skeleton only contain spatial information of actions, However, many actions are difficult to recognize from the spatial information





**Fig. 3.** The prediction scores of FGCN-spatial and FGCN-motion are fused for final action prediction.

alone, for example “wear a shoe” versus “take off a shoe”, “wear on glasses” versus “take off glasses” and *etc.* Inspired by [43], we model the spatial-temporal features not only exploiting spatial information but the temporal movement information of skeleton sequences. As in Section 3.1, the joint and bone of skeleton are denoted as a vector of coordinates. The movement of a joint or bone is defined as the difference of the vectors for the same joint or bone in consecutive frames along the temporal dimension.

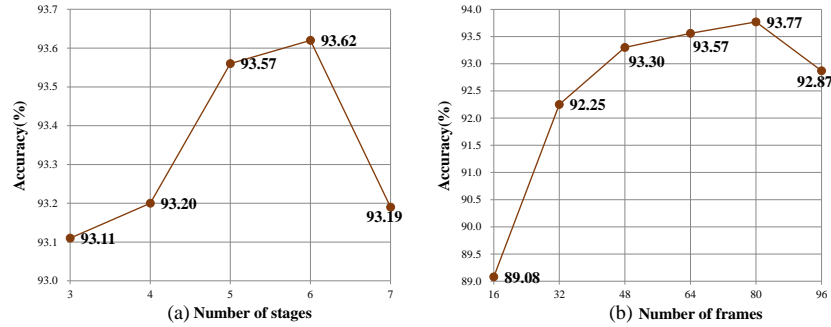
Given the joints and bones from two consecutive frames, denoted as  $v_{ti}$ ,  $v_{(t+1)i}$  and  $e_{v_{ti},v_{tj}}$ ,  $e_{v_{(t+1)i},v_{(t+1)j}}$  respectively, the movement of joints is defined as  $mv_{ti} = v_{(t+1)i} - v_{ti}$ . Similarly, the movement of bones is defined as  $me_t^{ij} = e_{v_{(t+1)i},v_{(t+1)j}} - e_{v_{ti},v_{tj}}$ . As the spatial information modeling, the motion information is formulated as a sequence of graphs  $S^m = \{G_1^m, G_2^m, \dots, G_{len}^m\}$ , where  $G_t^m = \{\mathbf{V}_t^m, \mathbf{E}_t^m\}$ ,  $\mathbf{V}_t^m = \{mv_{ti}\}_{i=1}^N$  and  $\mathbf{E}_t^m = \{me_t^{ij}\}_{(i,j) \in Q}$ . In this paper, the spatial graph  $S$  and the motion graph  $S^m$  are fed into two separate FGCN models to predict action labels. The model fed with spatial graphs  $S$  is denoted as FGCN-spatial, the other fed with temporal graphs  $S^m$  is denoted as FGCN-motion. The two models are finally fused by weighting the output scores of the softmax layers, as shown in Fig. 3.

## 4 Experiments

In this section, we evaluate the proposed FGCN method by conducting extensive experiments on three 3D skeleton action datasets, NTU-RGB+D, NTU-RGB+D120, and Northwestern-UCLA.

### 4.1 Datasets

**NTU-RGB+D** [41] is a widely used dataset for skeleton-based action recognition. The dataset contains more than 56,000 skeleton sequences categorized into 60 action classes. It provides 25 major body joints with 3D coordinates for every human in each frame. Two benchmark evaluations are recommended: cross-subject and cross-view. For cross-subject, both training and test sets consist of



**Fig. 4.** Evaluating the influence of two key factors on NTU-RGB+D, (a) influence of the number of stages, (b) influence of frame length in each stage.

20 subjects, and have 40,320 and 16,560 sequences respectively. The cross-view setup divides the data according to camera views. The training set has 37,920 sequences captured from the front and two side views, while the test set has 18,960 sequences captured from left and right 45 degree views.

**NTU-RGB+D120** [28] is currently the largest in-door captured 3D skeleton dataset. It is an extension of NTU-RGB+D with 120 action classes and more than 114,000 video samples. The newly added action classes make the action recognition more challenging. For example, different actions may have similar body motions but different subjects. There may be fine-grained hand or finger motions and so on. The dataset has 106 subjects and 32 setup IDs. Cross-subject and cross-setup benchmarks are defined. For cross-subject, 53 subjects constitute the training set, and the remaining 53 subjects constitute the test set. Analogously, the 32 setup IDs are also divided equally into two parts for training and testing in cross-setup.

**Northwestern-UCLA** [55] is a multi-view 3D event dataset captured simultaneously by three Kinect cameras from different viewpoints. This dataset includes 1494 video sequences covering 10 action categories performed by 10 subjects from 1 to 6 times. It provides 3D spatial coordinates of 20 major body joints. As reported in [55], we pick all samples from the first two cameras for training. The samples from the remaining cameras are for testing.

## 4.2 Implementation Details

All experiments are implemented with PyTorch deep learning framework. A stochastic gradient descent (SGD) optimizer is used during training with the batch size as 32, the momentum as 0.9, and the initial learning rate as 0.1. The learning rate is divided by 10 at the 40th and 60th epoch. The training process ends at the 80th epoch. In our experiments, the input video is divided into five stages temporally and 64 consecutive frames are sampled randomly from each stage to form an input clip. Ten graph convolutional layers are stacked at the front of the feedback block FGCB and these layers have the same configuration as the graph convolutional layers in ST-GCN [56]. The FGCB has four graph

convolutional layers (*i.e.*  $L = 4$ ). The spatial temporal kernel sizes and output channels of them are set as  $k_s = 3$ ,  $k_t = 3$  and  $m = 256$  respectively.

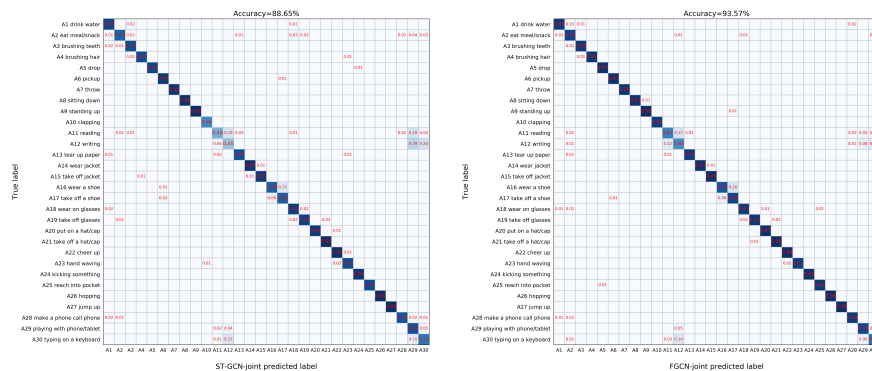
**Table 1.** Evaluating different temporal fusion strategies on NTU-RGB+D.

| Temporal Fusion Strategies | Weights |       |       |       |       | Cross-view(%) |
|----------------------------|---------|-------|-------|-------|-------|---------------|
|                            | $w_1$   | $w_2$ | $w_3$ | $w_4$ | $w_5$ |               |
| Last-win-all fusion        | 0       | 0     | 0     | 0     | 1     | 89.88         |
| Weight fusion-1            | 0.05    | 0.05  | 0.1   | 0.2   | 0.6   | 93.09         |
| Weight fusion-2            | 0.1     | 0.15  | 0.2   | 0.25  | 0.3   | 93.05         |
| Average fusion             | 0.2     | 0.2   | 0.2   | 0.2   | 0.2   | <b>93.57</b>  |

### 4.3 Ablation Study

In this section, we design four ablation experiments to evaluate the influence of different hyper-parameters, architecture and inputs on the performance of our FGCN model. These ablation experiments are all conducted on the challenging skeleton dataset NTU-RGB+D.

In the first experiment, we evaluate the influence of two key hyper-parameters on the performance of our FGCN model, *i.e.*, the number of stages and the length of the input clip in each stage. In Fig. 4(a), the performances of FGCN with different numbers of temporal stages are reported. The FGCN model achieves the best performance when the input video is divided into 6 stages with equal duration. In the subsequent experiments, we set the number of temporal stages at 5, to balance performance against computational cost. Similar performances are obtained with 6 temporal stages. In Fig. 4(b), we evaluate the performance of FGCN fed with different numbers of frames at each stage. Based on the similar model selection strategy in the last experiment, we set the frame length as 64 in the subsequent experiments to balance performance against computational cost.



**Fig. 5.** The confusion matrices of ST-GCN-joint and FGCN-joint on NTU-RGB+D.

**Table 2.** Evaluating the effectiveness of the FGCN model fed with different inputs on the NTU-RGB+D dataset.

| Models                   | Cross-subject(%) | Cross-view(%) |
|--------------------------|------------------|---------------|
| ST-GCN-joint [56]        | 81.5             | 88.3          |
| FGCN-joint               | 87.04            | 93.57         |
| FGCN-bone                | 86.96            | 93.22         |
| FGCN-joint+FGCN-bone     | 89.24            | 95.28         |
| FGCN-spatial             | 88.32            | 94.82         |
| FGCN-motion              | 85.96            | 93.57         |
| FGCN-spatial+FGCN-motion | <b>90.22</b>     | <b>96.25</b>  |

In the second experiment, we evaluate the effectiveness of different temporal fusion strategies in the FGCN model, *i.e.* last-win-all fusion, average fusion, and weight fusion. The experiment results are listed in Tab. 1. Among these three fusion strategies, the average fusion strategy achieves the best performance. Based on the results, we use the average fusion strategy to fuse the local predictions for the video-level prediction in the subsequent experiments.

In the third experiment, we evaluate the effectiveness of the proposed FGCN model fed with joints and bones. We first compare the proposed FGCN model with its baseline ST-GCN model. The two models have the same architecture and configuration of convolutional layers. As shown in the upper part of Tab. 2, the FGCN model fed with joint sequences of skeletons (FGCN-joint) outperforms the baseline model ST-GCN-joint by 5.54% and 5.27% on the cross-subject and cross-view benchmarks respectively. The confusion matrices for the former 30 classes are shown in Fig. 5, and the complete confusion matrices are shown in the supplementary materials. The improvements indicate that introducing feedback mechanism into GCNs is very effective for action recognition. Moreover, we fuse the softmax scores of two FGCN models, where one model is FGCN-joint, the other is FGCN-bone which is fed with the bone sequences. The fusion model FGCN-joint+FGCN-bone achieves a clear improvement, compared with FGCN-joint and FGCN-bone.

In the fourth experiment, we evaluate the effectiveness of FGCN model fed with the spatial information and the motion information, *i.e.* FGCN-spatial and FGCN-motion, on the NTU-RGB+D dataset. The experiment results of these two models and their fusion are reported in the under part of Tab. 2. Firstly, the FGCN-spatial model fed with spatial information (joints and bones) achieves 88.32% on cross-subject and 94.82% on cross-view. It is comparable with the performance of the FGCN-joint+FGCN-bone model that fuses the softmax scores of two models. Then, the FGCN-motion fed with the movement of joints and bones achieves 85.96% on cross-subject and 93.57% on cross-view. Finally, we fuse the softmax scores of FGCN-spatial and FGCN-motion. The FGCN-spatial+FGCN-motion achieves 90.22% on cross-subject and 96.25% on cross-view, and it achieves a clear improvement from both of FGCN-spatial and FGCN-motion.

**Table 3.** Comparisons with the state-of-the-art methods on NTU-RGB+D.

| Models                            | Cross-subject(%) | Cross-view(%) |
|-----------------------------------|------------------|---------------|
| ResNet152-3S (ICMEW 2017) [22]    | 85.0             | 92.3          |
| ST-GCN (AAAI 2018) [56]           | 81.5             | 88.3          |
| DPRL+GCNN (CVPR 2018) [50]        | 83.5             | 89.8          |
| SR-TSL (ECCV 2018) [46]           | 84.8             | 92.4          |
| PB-GCN (BMVC 2018) [51]           | 87.5             | 93.2          |
| Bayesian GC-LSTM (ICCV 2019) [61] | 81.8             | 89.0          |
| AS-GCN (CVPR 2019) [23]           | 86.8             | 94.2          |
| AGC-LSTM (CVPR 2019) [45]         | 89.2             | 95.0          |
| 2s-AGCN (CVPR 2019) [44]          | 88.5             | 95.1          |
| DGNN (CVPR 2019) [43]             | 89.9             | 96.1          |
| FGCN (ours)                       | <b>90.2</b>      | <b>96.3</b>   |

#### 4.4 Comparison with State-of-the-art

In this section, we compare the performance of the FGCN model with the recent state-of-the-art methods on the NTU-RGB+D dataset, the NTU-RGB+D120 dataset, and the Northwestern-UCLA dataset.

For the NTU-RGB+D dataset, we display the accuracy of skeleton based action recognition methods, such as CNN-based methods [22], RNN-based methods [46,45,61] and GCN based methods [23,43,44,56]. As shown in Tab. 3, the proposed FGCN model achieves 8.7% and 8.0% improvements on the cross-subject and cross-view benchmarks respectively over the most comparable method ST-GCN [56]. These improvements show the effectiveness of the proposed feedback framework in action recognition. Moreover, the FGCN model outperforms other recent state-of-the-art methods, such as AS-GCN [45], 2s-AGCN [44], and DGNN [43]. Our FGCN model achieves state-of-the-art performance on both cross-subject and cross-view benchmarks of the NTU-RGB+D dataset.

**Table 4.** Comparisons with the state-of-the-art methods on NTU-RGB+D120.

| Models                                       | Cross-subject(%) | Cross-setup(%) |
|--|------------------|----------------|
| Internal Feature Fusion (T-PAMI 2017) [30]   | 58.2             | 60.9           |
| Multi-Task Learning Network (CVPR 2017) [17] | 58.4             | 57.9           |
| Skeleton Visualization (PR 2017) [32]        | 60.3             | 63.2           |
| Two-Stream Attention LSTM (TIP 2017) [31]    | 61.2             | 63.3           |
| Multi-Task CNN with RotClips (TIP 2018) [18] | 62.2             | 61.8           |
| ST-GCN (AAAI 2018) (reported in [36])        | 72.4             | 71.3           |
| AS-GCN (CVPR 2019) (reported in [36])        | 77.7             | 78.9           |
| FSNet (T-PAMI 2019) [29]                     | 59.9             | 62.4           |
| TSRJI (SIBGRAPI 2019) [3]                    | 67.9             | 62.8           |
| LSTM-IRN (arXiv 2019) [38]                   | 77.7             | 79.6           |
| GVFE + AS-GCN (arXiv 2019) [36]              | 78.3             | 79.8           |
| FGCN (ours)                                  | <b>85.4</b>      | <b>87.4</b>    |

**Table 5.** Comparisons with the state-of-the-art methods on Northwestern-UCLA.

| Models                                | Accuracy(%) |
|---------------------------------------|-------------|
| Actionlet ensemble (T-PAMI 2013) [54] | 76.0        |
| Lie group (CVPR 2014 ) [53]           | 74.2        |
| HBRNN-L(CVPR 2015) [6]                | 78.5        |
| Skeleton Visualization (PR 2017) [32] | 86.1        |
| Ensemble TS-LSTM (ICCV 2017) [21]     | 89.2        |
| AGC-LSTM (CVPR 2019) [45]             | 93.3        |
| JS+JM+BS+BM (ICME 2019) [25]          | 91.3        |
| HiGCN (ICIG 2019) [14]                | 88.9        |
| MSNN (CSVT 2020) [42]                 | 89.4        |
| FGCN (ours)                           | <b>95.3</b> |

For the UTU-RGB+D120 dataset, the results on cross-subject and cross-setup benchmarks of the recent state-of-the-art methods are listed in Tab. 4. The proposed FGCN model achieves 85.4% on cross-subject and 87.4% on cross-setup and it outperforms the most comparable ST-GCN model [56] by 13.0% and 16.1% on the cross-subject and cross-setup benchmarks respectively. The FGCN model outperforms other state-of-the-art methods with much larger margins. For example, the FGCN model outperforms Two-Stream Attention LSTM [31] by over 24% on both cross-subject and cross-setup benchmarks, and outperforms the most recent work FSNet [29] by over 25% on both of cross-subject and cross-setup benchmarks.

For the typical 3D action recognition dataset Northwestern-UCLA, we compare the proposed FGCN model with the state-of-the-art methods in recent years. The results of these models are reported in Tab. 5. The FGCN model outperforms the part-based hierarchical recurrent neural network HBRNN-L [6] by 16.8%. The recent method AGC-LSTM proposes an attention enhanced graph convolutional LSTM network to capture discriminative features from the co-occurrence relationship between spatial configuration and temporal dynamics. The FGCN model outperforms it by 2%. Moreover, the FGCN model outperforms the most recent methods, such as JS+JM+BS+BM [25], HiGCN [14] and MSNN [42]. The proposed FGCN model achieves state-of-the-art performance on the Northwestern-UCLA dataset.

## 5 Conclusion

In this paper, we propose a novel FGCN model to extract effective spatial-temporal features of actions in a coarse-to-fine progressive process. Firstly, we propose a multi-stage temporal sampling strategy to sample sparse skeleton clips in multiple temporal stages and exploit graph convolutional layers to extract local spatial-temporal features for each stage. Then, we introduce the feedback mechanism into conventional GCNs by proposing the FGCB which is a local graph convolutional dense network. The FGCB transmits the semantic information from high-level layers to low-level layers and from the former stages

to the later stages. Moreover, the FGCN provides early predictions which help agents in many applications to make timely decisions on-the-fly. The proposed FGCN model is extensively evaluated on the NTU-RGB+D, NTU-RGB+D120 and Northwestern-UCLA datasets, indicating that the FGCN is effective for action recognition. It has achieved state-of-the-art performance on the three datasets.

## References

1. Ashford, S.J., Cummings, L.L.: Feedback as an individual resource: Personal strategies of creating information. *Organizational Behavior and Human Performance* **32**(3), 370–398 (1983)
2. Bruna, J., Zaremba, W., Szlam, A., Lecun, Y.: Spectral networks and locally connected networks on graphs. In: *International Conference on Learning Representations* (2014)
3. Caetano, C., Brémond, F., Schwartz, W.R.: Skeleton image representation for 3D action recognition based on tree structure and reference joints. In: *SIBGRAPI Conference on Graphics, Patterns and Images*. pp. 16–23. IEEE (2019)
4. Cao, Z., Simon, T., Wei, S.E., Sheikh, Y.: Realtime multi-person 2D pose estimation using part affinity fields. In: *IEEE Conference on Computer Vision and Pattern Recognition*. pp. 7291–7299 (2017)
5. Carreira, J., Agrawal, P., Fragkiadaki, K., Malik, J.: Human pose estimation with iterative error feedback. In: *IEEE Conference on Computer Vision and Pattern Recognition*. pp. 4733–4742 (2016)
6. Du, Y., Wang, W., Wang, L.: Hierarchical recurrent neural network for skeleton based action recognition. In: *IEEE Conference on Computer Vision and Pattern Recognition*. pp. 1110–1118 (2015)
7. Duvenaud, D.K., Maclaurin, D., Iparraguirre, J., Bombarell, R., Hirzel, T., Aspuru-Guzik, A., Adams, R.P.: Convolutional networks on graphs for learning molecular fingerprints. In: *Advances in Neural Information Processing Systems*. pp. 2224–2232 (2015)
8. Evangelidis, G., Singh, G., Horaud, R.: Skeletal quads: Human action recognition using joint quadruples. In: *IEEE International Conference on Pattern Recognition*. pp. 4513–4518 (2014)
9. Gilbert, C.D., Sigman, M.: Brain states: top-down influences in sensory processing. *Neuron* **54**(5), 677–696 (2007)
10. Han, W., Chang, S., Liu, D., Yu, M., Witbrock, M., Huang, T.S.: Image super-resolution via dual-state recurrent networks. In: *IEEE Conference on Computer Vision and Pattern Recognition*. pp. 1654–1663 (2018)
11. Haris, M., Shakhnarovich, G., Ukita, N.: Deep back-projection networks for super-resolution. In: *IEEE Conference on Computer Vision and Pattern Recognition*. pp. 1664–1673 (2018)
12. Henaff, M., Bruna, J., Lecun, Y.: Deep convolutional networks on graph-structured data. *Computer Science* (2015)
13. Hu, J.F., Zheng, W.S., Lai, J., Zhang, J.: Jointly learning heterogeneous features for RGB-D activity recognition. In: *IEEE Conference on Computer Vision and Pattern Recognition*. pp. 5344–5352 (2015)
14. Huang, L., Huang, Y., Ouyang, W., Wang, L.: Hierarchical graph convolutional network for skeleton-based action recognition. In: *Springer International Conference on Image and Graphics*. pp. 93–102 (2019)

15. Hupé, J., James, A., Payne, B., Lomber, S., Girard, P., Bullier, J.: Cortical feedback improves discrimination between figure and background by v1, v2 and v3 neurons. *Nature* **394**(6695), 784–787 (1998)
16. Karpathy, A., Toderici, G., Shetty, S., Leung, T., Sukthankar, R., Fei-Fei, L.: Large-scale video classification with convolutional neural networks. In: *IEEE Conference on Computer Vision and Pattern Recognition*. pp. 1725–1732 (2014)
17. Ke, Q., Bennamoun, M., An, S., Sohel, F., Boussaid, F.: A new representation of skeleton sequences for 3D action recognition. In: *IEEE Conference on Computer Vision and Pattern Recognition*. pp. 3288–3297 (2017)
18. Ke, Q., Bennamoun, M., An, S., Sohel, F., Boussaid, F.: Learning clip representations for skeleton-based 3D action recognition. *IEEE Transactions on Image Processing* **27**(6), 2842–2855 (2018)
19. Kipf, T.N., Welling, M.: Semi-supervised classification with graph convolutional networks. In: *International Conference on Learning Representations* (2017)
20. Lee, E.B., Markus, L.: Foundations of optimal control theory. Tech. rep., Minnesota Univ Minneapolis Center For Control Sciences (1967)
21. Lee, I., Kim, D., Kang, S., Lee, S.: Ensemble deep learning for skeleton-based action recognition using temporal sliding LSTM networks. In: *IEEE International Conference on Computer Vision*. pp. 1012–1020 (2017)
22. Li, B., Dai, Y., Cheng, X., Chen, H., Lin, Y., He, M.: Skeleton based action recognition using translation-scale invariant image mapping and multi-scale deep CNN. In: *IEEE International Conference on Multimedia and Expo Workshops*. pp. 601–604 (2017)
23. Li, M., Chen, S., Chen, X., Zhang, Y., Wang, Y., Tian, Q.: Actional-structural graph convolutional networks for skeleton-based action recognition. In: *IEEE Conference on Computer Vision and Pattern Recognition*. pp. 3595–3603 (2019)
24. Li, S., Li, W., Cook, C., Zhu, C., Gao, Y.: Independently recurrent neural network (indrnn): Building a longer and deeper RNN. In: *IEEE Conference on Computer Vision and Pattern Recognition*. pp. 5457–5466 (2018)
25. Li, Y., Xia, R., Liu, X., Huang, Q.: Learning shape-motion representations from geometric algebra spatio-temporal model for skeleton-based action recognition. In: *IEEE International Conference on Multimedia and Expo*. pp. 1066–1071 (2019)
26. Li, Z., Yang, J., Liu, Z., Yang, X., Jeon, G., Wu, W.: Feedback network for image super-resolution. In: *IEEE Conference on Computer Vision and Pattern Recognition*. pp. 3867–3876 (2019)
27. Liu, H., Tu, J., Liu, M.: Two-stream 3D convolutional neural network for skeleton-based action recognition. *arXiv preprint arXiv:1705.08106* (2017)
28. Liu, J., Shahroudy, A., Perez, M.L., Wang, G., Duan, L.Y., Chichung, A.K.: NTU RGB+D 120: A large-scale benchmark for 3D human activity understanding. *IEEE Transactions on Pattern Analysis and Machine Intelligence* (2019)
29. Liu, J., Shahroudy, A., Wang, G., Duan, L.Y., Chichung, A.K.: Skeleton-based on-line action prediction using scale selection network. *IEEE Transactions on Pattern Analysis and Machine Intelligence* (2019)
30. Liu, J., Shahroudy, A., Xu, D., Kot, A.C., Wang, G.: Skeleton-based action recognition using spatio-temporal LSTM network with trust gates. *IEEE Transactions on Pattern Analysis and Machine Intelligence* **40**(12), 3007–3021 (2017)
31. Liu, J., Wang, G., Duan, L.Y., Abdiyeva, K., Kot, A.C.: Skeleton-based human action recognition with global context-aware attention LSTM networks. *IEEE Transactions on Image Processing* **27**(4), 1586–1599 (2017)
32. Liu, M., Liu, H., Chen, C.: Enhanced skeleton visualization for view invariant human action recognition. *Elsevier Pattern Recognition* **68**, 346–362 (2017)



33. Luo, J., Wang, W., Qi, H.: Group sparsity and geometry constrained dictionary learning for action recognition from depth maps. In: IEEE International Conference on Computer Vision. pp. 1809–1816 (2013)
34. Niepert, M., Ahmed, M., Kutzkov, K.: Learning convolutional neural networks for graphs. In: IEEE International Conference on Machine Learning. pp. 2014–2023 (2016)
35. Ohn-Bar, E., Trivedi, M.: Joint angles similarities and HOG2 for action recognition. In: IEEE Conference on Computer Vision and Pattern Recognition Workshops. pp. 465–470 (2013)
36. Papadopoulos, K., Ghorbel, E., Aouada, D., Ottersten, B.: Vertex feature encoding and hierarchical temporal modeling in a spatial-temporal graph convolutional network for action recognition. arXiv preprint arXiv:1912.09745 (2019)
37. Parlos, A.G., Chong, K.T., Atiya, A.F.: Application of the recurrent multilayer perceptron in modeling complex process dynamics. IEEE Transactions on Neural Networks **5**(2), 255–266 (1994)
38. Perez, M., Liu, J., Kot, A.C.: Interaction relational network for mutual action recognition. arXiv preprint arXiv:1910.04963 (2019)
39. Rahmani, H., Mahmood, A., Huynh, D.Q., Mian, A.: Real time action recognition using histograms of depth gradients and random decision forests. In: IEEE Winter Conference on Applications of Computer Vision. pp. 626–633 (2014)
40. Rahmani, H., Mian, A.: Learning a non-linear knowledge transfer model for cross-view action recognition. In: IEEE Conference on Computer Vision and Pattern Recognition. pp. 2458–2466 (2015)
41. Shahroudy, A., Liu, J., Ng, T.T., Wang, G.: NTU RGB+D: A large scale dataset for 3D human activity analysis. In: IEEE Conference on Computer Vision and Pattern Recognition. pp. 1010–1019 (2016)
42. Shao, Z., Li, Y., Zhang, H.: Learning representations from skeletal self-similarities for cross-view action recognition. IEEE Transactions on Circuits and Systems for Video Technology (2020)
43. Shi, L., Zhang, Y., Cheng, J., Lu, H.: Skeleton-based action recognition with directed graph neural networks. In: IEEE Conference on Computer Vision and Pattern Recognition. pp. 7912–7921 (2019)
44. Shi, L., Zhang, Y., Cheng, J., Lu, H.: Two-stream adaptive graph convolutional networks for skeleton-based action recognition. In: IEEE Conference on Computer Vision and Pattern Recognition. pp. 12026–12035 (2019)
45. Si, C., Chen, W., Wang, W., Wang, L., Tan, T.: An attention enhanced graph convolutional LSTM network for skeleton-based action recognition. In: IEEE Conference on Computer Vision and Pattern Recognition. pp. 1227–1236 (2019)
46. Si, C., Jing, Y., Wang, W., Wang, L., Tan, T.: Skeleton-based action recognition with spatial reasoning and temporal stack learning. In: Springer European Conference on Computer Vision. pp. 103–118 (2018)
47. Simonyan, K., Zisserman, A.: Two-stream convolutional networks for action recognition in videos. In: Advances in Neural Information Processing Systems. pp. 568–576 (2014)
48. Song, S., Lan, C., Xing, J., Zeng, W., Liu, J.: An end-to-end spatio-temporal attention model for human action recognition from skeleton data. In: AAAI Conference on Artificial Intelligence. pp. 4263–4270 (2017)
49. Stollenga, M.F., Masci, J., Gomez, F., Schmidhuber, J.: Deep networks with internal selective attention through feedback connections. In: Advances in Neural Information Processing Systems. pp. 3545–3553 (2014)

50. Tang, Y., Tian, Y., Lu, J., Li, P., Zhou, J.: Deep progressive reinforcement learning for skeleton-based action recognition. In: IEEE Conference on Computer Vision and Pattern Recognition. pp. 5323–5332 (2018)
51. Thakkar, K., Narayanan, P.: Part-based graph convolutional network for action recognition. In: British Machine Vision Conference. pp. 1–13 (2018)
52. Toshev, A., Szegedy, C.: Deeppose: Human pose estimation via deep neural networks. In: IEEE Conference on Computer Vision and Pattern Recognition. pp. 1653–1660 (2014)
53. Vemulapalli, R., Arrate, F., Chellappa, R.: Human action recognition by representing 3D skeletons as points in a lie group. In: IEEE Conference on Computer Vision and Pattern Recognition. pp. 588–595 (2014)
54. Wang, J., Liu, Z., Wu, Y., Yuan, J.: Learning actionlet ensemble for 3D human action recognition. *IEEE Transactions on Pattern Analysis and Machine Intelligence* **36**(5), 914–927 (2013)
55. Wang, J., Nie, X., Xia, Y., Wu, Y., Zhu, S.C.: Cross-view action modeling, learning and recognition. In: IEEE Conference on Computer Vision and Pattern Recognition. pp. 2649–2656 (2014)
56. Yan, S., Xiong, Y., Lin, D.: Spatial temporal graph convolutional networks for skeleton-based action recognition. In: AAAI Conference on Artificial Intelligence. pp. 7444–7452 (2018)
57. Zamir, A.R., Wu, T.L., Sun, L., Shen, W.B., Shi, B.E., Malik, J., Savarese, S.: Feedback networks. In: IEEE Conference on Computer Vision and Pattern Recognition. pp. 1308–1317 (2017)
58. Zhang, P., Lan, C., Xing, J., Zeng, W., Xue, J., Zheng, N.: View adaptive recurrent neural networks for high performance human action recognition from skeleton data. In: IEEE International Conference on Computer Vision. pp. 2117–2126 (2017)
59. Zhang, Z.: Microsoft kinect sensor and its effect. *IEEE Multimedia* **19**(2), 4–10 (2012)
60. Zhao, R., Ali, H., Van der Smagt, P.: Two-stream RNN/CNN for action recognition in 3D videos. In: IEEE International Conference on Intelligent Robots and Systems. pp. 4260–4267 (2017)
61. Zhao, R., Wang, K., Su, H., Ji, Q.: Bayesian graph convolution LSTM for skeleton based action recognition. In: IEEE International Conference on Computer Vision. pp. 6882–6892 (2019)