

# 1 Attention module-based spatial temporal graph convolutional 2 networks for skeleton-based action recognition

3  
4 Yinghui Kong,<sup>a</sup> Li Li,<sup>a</sup> Ke Zhang,<sup>a</sup> Qiang Ni,<sup>b</sup> Jungong Han<sup>b</sup>

5 <sup>a</sup>North China Electric Power University, Department of Electronic and Communication Engineering, North  
6 Yonghua Road #619, Baoding, China, 071000

7 <sup>b</sup>Lancaster University, School of Computing and Communications, Bailrigg, Lancaster, United Kingdom, LA1 4YW  
8

9 **Abstract.** Skeleton-based action recognition is a significant direction of human action recognition, because the  
10 skeleton contains important information for recognizing action. The spatial temporal graph convolutional networks  
11 (ST-GCN) automatically learn both the temporal and spatial features from the skeleton data, and achieve remarkable  
12 performance for skeleton-based action recognition. However, ST-GCN just learn local information on a certain  
13 neighborhood, but does not capture the correlation information between all joints (i.e., global information).  
14 Therefore, we need to introduce global information into the spatial temporal graph convolutional networks. In this  
15 work, we propose a model of dynamic skeletons called attention module-based Spatial Temporal Graph  
16 Convolutional Networks (AM-STGCN), which solves these problems by adding attention module. The attention  
17 module can capture some global information, which brings stronger expressive power and generalization capability.  
18 Experimental results on two large-scale datasets, Kinetics and NTU-RGB+D, demonstrate that our model achieves  
19 significant improvements over previous representative methods.

20  
21 **Keywords:** action recognition, spatial temporal graph convolution network, non-local neural network, attention  
22 module.

23  
24 \*Yinghui Kong, E-mail: [kongyhbd2015@ncepu.edu.cn](mailto:kongyhbd2015@ncepu.edu.cn)  
25

## 26 27 1 Introduction

28 Action recognition technology plays an increasingly important role in many fields such as  
29 intelligent monitoring, human-computer interaction, video sequence understanding, and medical  
30 health. Video action recognition technology is challenged by factors such as occlusion, dynamic  
31 background, mobile camera, angle of view and illumination change.

32 Before the advent of deep learning, the best algorithm for human action recognition in video  
33 was iDT<sup>1,2</sup>, and the subsequent works were basically improved based on the iDT method. Human  
34 action recognition uses multiple modalities of data such as appearance, depth, optical flows, and  
35 body skeletons.<sup>3</sup> With the continuous development of deep learning and its excellent

36 performance in image understanding tasks, more and more researchers are beginning to use deep  
37 learning methods to solve the problem of video analysis. Action recognition methods based on  
38 RGB video or optical flows, such as Two-Stream<sup>4,5</sup>, C3D<sup>6</sup>, I3D<sup>7</sup>, RNN<sup>8</sup> methods, are greatly  
39 affected by illumination, scene and camera lens movement, so it is difficult to describe the  
40 motion of the human body in the sequence, the recognition performance in some complex  
41 datasets needs to be improved. In recent years, due to the cost reduction of depth sensors (such as  
42 Kinect) and the emergence of real-time human pose estimation algorithms, skeleton-based action  
43 recognition has become more and more popular.

44 Skeleton-based action recognition methods have been widely studied and paid attention due  
45 to its strong adaptability to dynamic environments and complex backgrounds. Traditional  
46 methods<sup>9,10</sup> require hand-crafted features and traversal rules, which are less efficient. Ordinary  
47 deep learning-based methods<sup>11-20</sup> manually structure the skeleton into joint coordinate vectors or  
48 pseudo-images, which are then sent to the RNN or CNN network for prediction of the action  
49 categories. The human skeleton is naturally constructed as a graph in a non-Euclidean space, in  
50 which the joint acts as a node, and the edge is constructed according to the natural connection  
51 relationship of the human body. Recently, the Graph Convolutional Networks (GCN) have  
52 extended convolution operations from images to graph structures, and have been successfully  
53 applied to many applications. For skeleton-based action recognition, GCN-based methods  
54 contain ST-GCN<sup>3</sup>, STGC<sup>21</sup>, SR-TSL<sup>22</sup>, AGCN<sup>23</sup>, PB-GCN<sup>24</sup>, GR-GCN<sup>25</sup> and DPRL+GCNN<sup>26</sup>.  
55 ST-GCN applied GCN for skeleton-based action recognition task and directly model the original  
56 skeleton data, it extended graph neural networks to a spatial-temporal graph model, and obtained  
57 better action representations. Compared to ordinary deep learning-based methods, GCN-based  
58 methods can better express the dependencies between joints. However, the convolution operation

59 in the ST-GCN method is performed on the 1-neighbor of the root node and cannot capture  
60 global information. For the action categories in which the interaction joints are not in the same  
61 neighborhood, such as brushing, clapping, but there are relations between these nonadjacent  
62 joints, attention mechanism can learn these relations. Paying more attention to those joints may  
63 improve recognition performance. [Attention modules that work well include non-local neural  
64 networks<sup>27</sup>, Interaction-aware attention<sup>28</sup>, CBAM<sup>29</sup>, SENet<sup>30</sup> etc.](#)

65 In order to solve this problem, we propose an improved method based on ST-GCN, which is  
66 attention module-based Spatial Temporal Graph Convolutional Networks (AM-STGCN).  
67 Attention module helps the model focus on all positions and learn different weights for each  
68 position. In AM-STGCN, we add the non-local neural network as an attention module after the  
69 convolution operation of the baseline model ST-GCN to learn the feature representation with  
70 long-range dependencies. In addition, we discussed the effects of adding attention blocks to  
71 different layers, as well as the effects of adding multiple attention blocks. We did a lot of  
72 experimentation and analysis, and finally got the best strategy. The experimental results on two  
73 large-scale action recognition datasets Kinetics<sup>31</sup> and NTU-RGB+D<sup>32</sup> show that AM-STGCN can  
74 significantly outperform ST-GCN in action recognition.

75 In the remainder of the paper, we first provide some related work in Sec. 2, and then  
76 introduce the original ST-GCN model and our AM-STGCN model in Sec. 3. We summarize and  
77 analyze the experimental results in Sec. 4. Finally, we draw conclusions and point out future  
78 research direction in Sec. 5.

## 79 2 Related Work

### 80 2.1 Action Recognition Based on RGB Video or Optical Flows

81 Most previous studies were based on RGB video or optical flows. Traditional action recognition  
82 methods are mostly based on optical flows, and the representative algorithm is iDT<sup>1,2</sup>. DT  
83 algorithm utilize optical flow field to obtain some trajectories in the video sequence, then extract  
84 the HOF, HOG, MBH and trajectory characteristics along the trajectory. IDT improves dense  
85 trajectories by explicitly estimating camera motion. Then, some methods based on deep learning  
86 gradually appeared, and their performance was much better than traditional methods. Two-  
87 stream method was originally proposed by Simonyan et al.<sup>4</sup>, and Feichtenhofer et al.<sup>5</sup> improved  
88 the model. Two-stream method utilizes both appearance and optical flows information: in spatial  
89 stream, in the form of appearance on a single frame, the scene and target information depicted by  
90 video are carried; in temporal stream, the motion of the observer (camera) and the target are  
91 expressed in the form of multi-frame optical flows. Tran et al.<sup>6</sup> adopted 3D convolution and 3D  
92 pooling to construct a network, which can directly process video, and its efficiency is much  
93 higher than other methods. Carreira et al.<sup>7</sup> proposed a model named “I3D” based on Inceptionv1,  
94 which inflates Inceptionv1’s filters and pooling kernels into 3D, leading to very deep, naturally  
95 spatiotemporal classifiers. Du et al.<sup>8</sup> introduced a novel pose-attention mechanism to adaptively  
96 learn pose-related features at every time-step action prediction of RNNs.

97 Although action recognition methods based on RGB video or optical flows perform high  
98 performance, there are still some problems. For example, it is susceptible to background,  
99 illumination and appearance changes, and extract optical flow information requires high  
100 computational cost.

## 101 2.2 *Skeleton-based Action Recognition*

102 The human skeleton can provide a very good representation of the human body motions, which  
103 is beneficial to the analysis of human actions. On the one hand, skeleton data is inherently robust  
104 in background noise, and provides abstract and high-level features of human motion. On the  
105 other hand, the size of the skeleton data is very small compared to RGB data, which allows us to  
106 design a lightweight and hardware-friendly model.

107 Skeleton-based action recognition approaches can be categorized into traditional methods  
108 and deep learning methods. Deep learning methods contain RNN based methods, CNN based  
109 methods and graph convolutional network (GCN) based methods.

110 Some traditional methods shown in Refs. 9 and 10 require hand-crafted features and traversal  
111 rules to achieve skeleton action recognition. With the development of deep learning, RNN based  
112 methods appears gradually. Du et al.<sup>11</sup> divided the human skeleton into five parts according to  
113 human physical structure, and then separately feeded them to five bidirectionally recurrently  
114 connected subnets. Song et al.<sup>12</sup> proposed an end-to-end spatial and temporal attention model,  
115 which learns to selectively focus on discriminative joints of skeleton within each frame of the  
116 inputs and pays different levels of attention to the outputs of different frames. Zhang et al.<sup>13</sup>  
117 designed a view adaptive recurrent neural network (RNN) with LSTM architecture, which  
118 enables the network itself to adapt to the most suitable observation viewpoints from end to end.  
119 In recent years, a number of CNN based approaches have also emerged. Kim et al.<sup>14</sup> re-designed  
120 the original TCN by factoring out the deeper layers into additive residual terms which yields  
121 both interpretable hidden representations and model parameters. Liu et al.<sup>15</sup> proposed an  
122 enhanced skeleton visualization method to represent a skeleton sequence as a series of visual and  
123 motion enhanced color images, which implicitly describe spatio-temporal skeleton joints in a

124 compact yet distinctive manner. Li et al.<sup>16</sup> designed a novel skeleton transformer module to  
125 rearrange and select important skeleton joints automatically. Li et al.<sup>17</sup> proposed an end-to-end  
126 convolutional co-occurrence feature learning framework to aggregate different levels of  
127 contextual information. Liu et al.<sup>18</sup> proposed a recurrent attention mechanism for their GCA-  
128 LSTM network, which is able to selectively focus on the informative joints in the action  
129 sequence with the assistance of global contextual information. Xie et al.<sup>19</sup> designed a temporal-  
130 then-spatial recalibration scheme, resulting in an end-to-end Memory Attention Networks  
131 (MANs) which consist of a Temporal Attention Recalibration Module (TARM) and a Spatio-  
132 Temporal Convolution Module (STCM). Zheng et al.<sup>20</sup> designed an adaptive attentional module  
133 to focus attention on the most discriminative parts in the single skeleton. Although RNN based  
134 methods has a strong ability to model sequence data, and CNN based methods has good  
135 parallelism and easier training process, however, neither CNN nor RNN fully represent the  
136 structure of the skeleton.

137 Recently, some methods based on graph convolution have appeared, and the effect has been  
138 improved obviously. Yan et al.<sup>3</sup> directly simulated the original skeleton using the graph  
139 convolution, which eliminates the need for manual part assignment, and it is easier to design and  
140 potent to learn better action representations. Li et al.<sup>21</sup> designed multi-scale convolutional filters  
141 to encode the graph structure data, and proposed a recursive graph convolution model. Si et al.<sup>22</sup>  
142 utilized a spatial reasoning network to capture the high-level spatial structural features within  
143 each frame, and utilized a composition of multiple skip-clip LSTMs to model the detailed  
144 temporal dynamics of skeleton sequences. In order to design individual graphs for different  
145 samples, Shi et al.<sup>23</sup> introduced non-local neural networks into graph convolution operation to  
146 model the multi-level semantic information, which brings more flexibility and generality.

147 Thakkar et al.<sup>24</sup> divided the skeleton graph into four subgraphs, and used relative coordinates and  
148 temporal displacements as features at each node instead of 3D joint coordinates which improves  
149 action recognition performance. Gao et al.<sup>25</sup> constructed a generalized graph via spectral graph  
150 theory to capture the space-time variation. Tang et al.<sup>26</sup> proposed a deep progressive  
151 reinforcement learning (DPRL) method to extract key frames, and employed the graph-based  
152 convolutional neural network to capture the dependency between the joints for action recognition.

### 153 **3 Methodology**

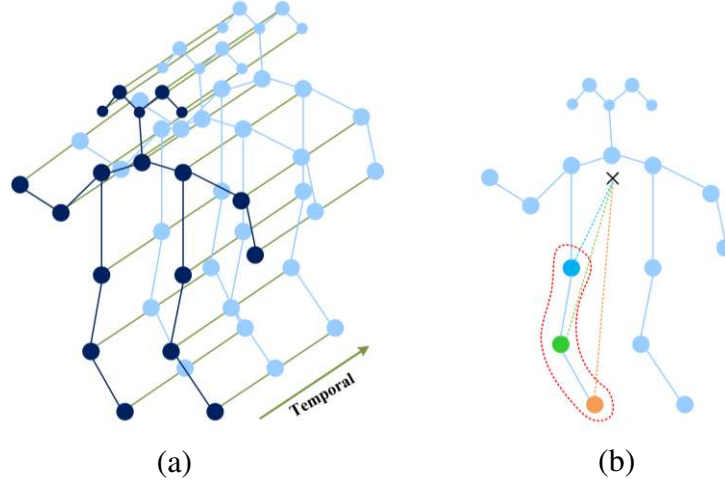
154 We briefly describe the original spatial temporal graph convolutional networks (ST-GCN) in Sec.  
155 3.1. And in Sec. 3.2, we give a briefly description about the methods of utilizing the attention  
156 module to boost the performance, and propose the improved model -- attention module-based  
157 spatial temporal graph convolution network (AM-STGCN).

#### 158 *3.1 Spatial-Temporal Graph Convolutional Networks (ST-GCN)*

159 As shown in Ref. 3, the authors take joints as nodes and the connections between nodes as edges  
160 to construct the skeleton graph. Fig. 1 (a) shows an example of a spatial-temporal skeleton graph.  
161 In one frame, the natural connections between the joints (i.e., the human bones) act as spatial  
162 edges; in adjacent frames, the same joints are joined as temporal edges. The property of each  
163 node is the coordinate vector of the joint. Multi-layers spatial-temporal graph convolution  
164 operation is applied to the spatial-temporal skeleton graph to obtain advanced feature map, and  
165 then use the SoftMax classifier to predict the action category.

166 ST-GCN applies the spatial configuration partitioning strategy shown in Fig. 1(b) in frame.  
167 The spatial configuration partitioning strategy divides the node's 1-neighbor into three subsets: 1)  
168 the root node (green dot); 2) the centripetal subset (blue dots): the neighboring nodes closer to

169 the gravity center of the skeleton (black cross); 3) the centrifugation subset (yellow dots): the  
 170 neighboring nodes that are further to the gravity center of the skeleton. Each color in the Fig. 1(b)  
 171 corresponds to a specific learnable weight vector. The authors of ST-GCN propose three  
 172 partitioning strategy, and it has been proved that the spatial configuration partitioning strategy  
 173 shown in Fig. 1(b) is the best, so this work directly adopts this strategy.



174  
 175  
 176 **Fig. 1** (a) Spatial temporal graph of the skeleton. (b) Partitioning strategy, different colors represent different  
 177 subsets.

178 Spatial graph convolution is formulated as:

$$179 \quad f_{out}(v_{ti}) = \sum_{v_{ij} \in B(v_{ti})} \frac{1}{Z_{ti}(v_{ij})} f_{in}(v_{ij}) \cdot w(l_{ti}(v_{ij})), \quad (1)$$

180 where  $f$  is the feature map.  $v_{ti}$  is the node of the graph.  $B(v_{ti})$  is the sampling area, which is  
 181 defined as the 1-neighbor set of joint nodes. The neighbor set  $B(v_{ti})$  of a joint node  $v_{ti}$  is  
 182 partitioned into a fixed number of  $K$  subsets, where each subset has a numeric label.<sup>3</sup> The  
 183 mapping function  $l_{ti}$  maps a node in the neighborhood to its subset label. The weight function  $w$   
 184 gives different weights according to different  $l_{ti}$  values. The normalizing term  $Z_i(v_j)$  equals the  
 185 cardinality of the corresponding subset.



186 To model the spatial temporal dynamics within skeleton sequence, since the number of  
 187 neighbors per node is fixed at 2 (the corresponding joint in the previous and subsequent frames),  
 188 it is directly to perform the graph convolution similar to the classical convolution operation,  
 189 concretely, we perform a  $K_c \times 1$  convolution on the output feature map computed above.<sup>23</sup>

190 In the single frame case, ST-GCN with the spatial configuration partitioning strategy can be  
 191 implemented with the following formula:

$$192 \quad f_{out} = \sum_j (\Lambda_j^{-\frac{1}{2}} A_j \Lambda_j^{-\frac{1}{2}}) \otimes M_j f_{in} W_j. \quad (2)$$

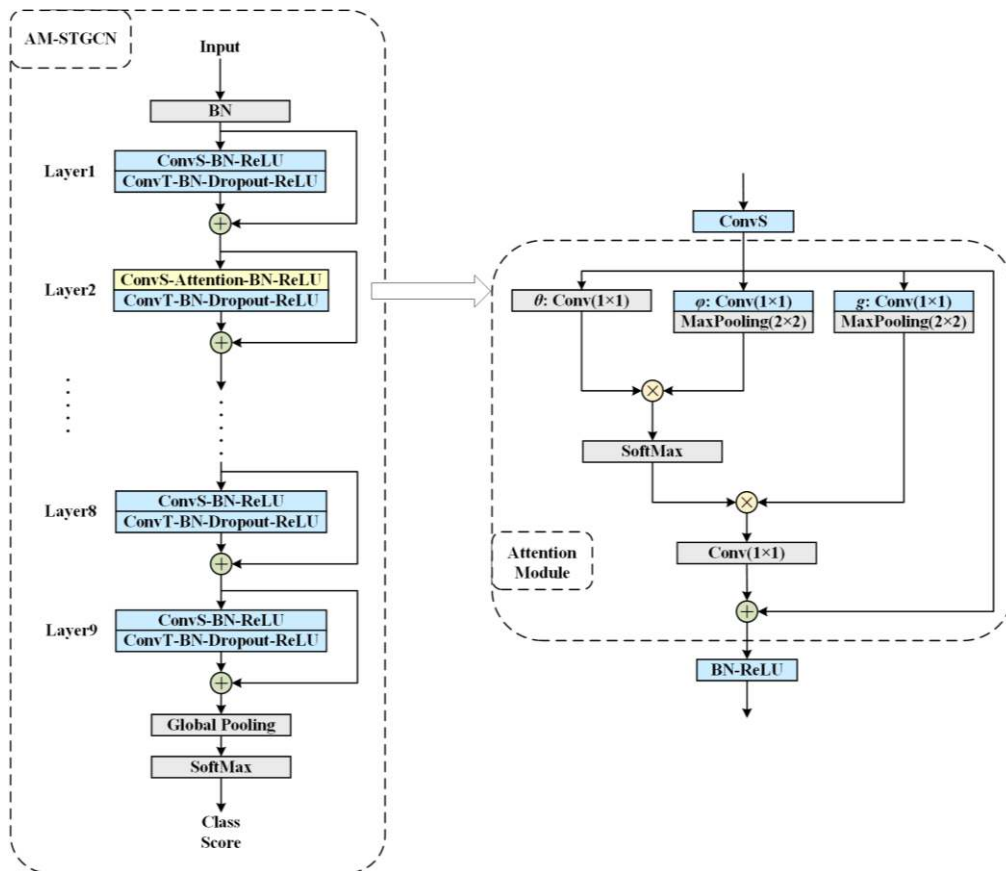
193 In formula 2,  $f$  is the  $C_{in} \times T \times V$  feature map where  $V$  denotes the number of nodes,  $T$  denotes the  
 194 temporal length and  $C_{in}$  denotes the number of input channels.  $A$  is the  $18 \times 18 \times 3$  adjacency  
 195 matrix, whose element  $A_{ij}$  indicates whether the node  $v_i$  is in the subset of node  $v_j$ .  $A_0 = I$   
 196 denotes the self-connections of vertexes,  $A_1$  denotes the connections of centripetal subset  
 197 and  $A_2$  denotes the centrifugal subset.  $\Lambda_j^{ii} = \sum_k (A_j^{ki}) + \alpha$  is the normalized diagonal matrix,  $\alpha$  is  
 198 set to 0.001 to avoid the empty rows in  $A$ .  $W_j$  is the  $C_{out} \times C_{in} \times 1 \times 1$  weight vector of  
 199 the  $1 \times 1$  convolution operation.  $M$  is a  $V \times V$  learnable attention map which indicates the  
 200 importance of each node.  $\otimes$  denotes the element-wise product between two matrixes. This  
 201 means that if one of the elements in  $A$  is 0, then whatever the value of  $M$  is, it will always be 0.  
 202 So  $M$  just operates in the 1-neighbor of the root node.

### 203 3.2 Attention Module-based Spatial Temporal Graph Convolution Network

204 In the spatial temporal graph convolution model, the receptive field of the convolution operation  
 205 is the 1-neighbor of the root node, so it only captures local features. However, in different  
 206 sample of different action classes, the relationship between the joints is not limited to the 1-

207 neighbor of the joint. For example, for many actions such as combing hair, brushing teeth, the  
 208 relationship between the hand and the head may be important. In order to solve this problem, we  
 209 introduce the idea of non-local neural network<sup>27</sup>, make some improvements to the ST-GCN  
 210 model, and then propose AM-STGCN skeleton-based action recognition method based on the  
 211 non-local attention mechanism, which directly focuses on the features of all joints, and get more  
 212 efficient features by attention operations.

213  
 214



215  
 216

**Fig. 2** The structure of AM-STGCN.

217 Fig. 2 shows the network structure of AM-STGCN, where we add the attention module after  
 218 the spatial convolution operation (ConvS) of Layer2. The model consists of nine layers of spatial  
 219 temporal graph convolution operators. The first three layers have 64 output channels, the middle  
 220 three layers have 128 output channels, and the last three layers have 256 output channels. Each

221 layer of AM-STGCN includes the spatial convolution operation (ConvS) and the temporal  
222 convolution operation (ConvT). The residual connection<sup>33</sup> is added on each layer.

223 Non-local neural network is a versatile, flexible building block, it can be easily embedded  
224 into existing 2D and 3D convolutional networks to improve or visualize related CV tasks. This  
225 allows us to combine global and local information to build richer hierarchy. In Fig. 2, the right  
226 side is our attention module, which is used to capture the correlation between all joints. We  
227 construct the attention module mainly following the idea of non-local neural network: first, linear  
228 mapping is conducted on the feature map of ConvS, which is implemented as  $1 \times 1$  convolution,  
229 and then get the  $\theta$ ,  $\phi$ ,  $g$  features; second, we perform a matrix point multiplication operation on  
230  $\theta$  and  $\phi$  to calculate the autocorrelation in the feature, and then carry out Softmax operation to  
231 obtain the self-attention coefficient; third, the attention coefficient is multiplied back into the  
232 feature matrix  $g$ ; at last, residual connection is established with the original input feature map,  
233 and then we get a new set of features. Specifically, we add  $2 \times 2$  MaxPooling operation after  $\theta$ ,  
234  $\phi$  features to reduce computational cost. Such an attention module is called one attention block,  
235 and multiple attention blocks will be used in the work. How many attention blocks are added to  
236 the model and where they are added will be analyzed in detail in Sec. 4, and the experimental  
237 results are given at the same time.

## 238 4 Experiments and Analysis

239 In this section, we evaluate the performance of the AM-STGCN model. In order to compare with  
240 the baseline model ST-GCN, our experiments are performed on the same two large-scale action  
241 recognition datasets: the human action dataset Kinetics<sup>31</sup> is the largest unconstrained action  
242 recognition dataset up to now, and NTU-RGB+D<sup>32</sup> is the largest constrained indoor captured

243 action recognition dataset. First, we conduct a detailed ablation study of the Kinetics dataset to  
244 analyze the contribution of the proposed model to recognition performance. Then, the  
245 corresponding experiments are carried out on the NTU-RGB+D dataset to verify whether the  
246 proposed model has certain generalization ability. Finally, we compare AM-STGCN with ST-  
247 GCN and some state-of-the-art results of skeleton-based action recognition on Kinetics and  
248 NTU-RGB+D. All experiments were performed on PyTorch deep learning framework using two  
249 1080Ti GPUs.

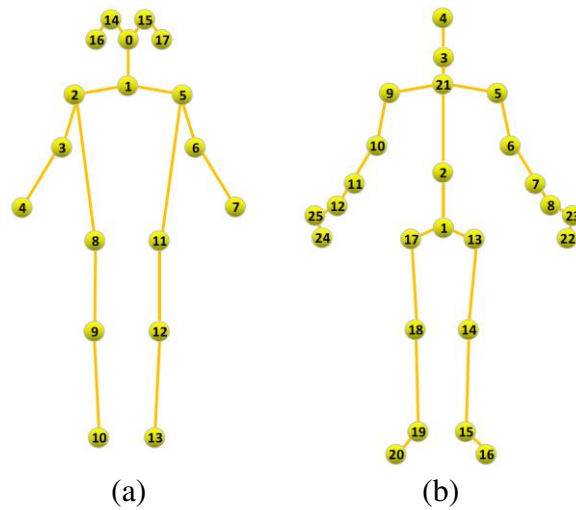
#### 250 4.1 Datasets

251 **Kinetics<sup>31</sup>**: Kinetics is a large human action dataset that contains 400 action classes taken from  
252 different YouTube video, each class with at least 400 video clips, each clip lasts about 10  
253 seconds<sup>31</sup>. These actions include the interaction between people and objects, such as playing an  
254 instrument, and the interaction between people, such as shaking hands.

255 The Kinetics dataset only provides raw video clips and does not provide skeleton joint data.  
256 As shown in Ref. 3, they use the public available OpenPose<sup>34</sup> toolbox to estimate the location of  
257 18 joints on every frame of the clips. In this work, we use the Kinetics-skeleton dataset provided  
258 by the author of ST-GCN, which marks the position of 18 joints in each frame. The dataset  
259 provides a training set of 240,000 clips and a validation set of 20,000 clips. In accordance with  
260 the recommendations in Ref. 31, in this work, we train the model on the training set and report  
261 the top-1 and top-5 recognition accuracies on the validation set.

262 Fig. 3(a) shows the joint label of the Kinetics-skeleton dataset. The joint labels are: 0 nose, 1  
263 neck, 2 right shoulder, 3 right elbow, 4 right wrist, 5 left shoulder, 6 left elbow, 7 left wrist, 8  
264 right hip, 9 right knee, 10 right ankle, 11 left hip, 12 left knee, 13 left ankle, 14 right eye, 15 left  
265 eye, 16 right ear, 17 left ear.

266 **NTU-RGB+D<sup>32</sup>**: NTU-RGB+D is the largest dataset with 3D joint annotations currently used  
 267 for human action recognition tasks. The dataset contains 60 action classes with a total of 56,000  
 268 action clips. All of these clips are performed by 40 volunteers in a constrained lab environment,  
 269 and captured by 3 cameras of the same height but from different horizontal angles:  $-45^\circ$ ,  $0^\circ$ ,  
 270  $45^\circ$ <sup>32</sup>. The dataset provides the 3D joint location of each frame detected by the Kinect depth  
 271 sensor. There are 25 joints per subject in the skeleton sequence. Each clip is guaranteed to have a  
 272 maximum of 2 subjects.



273  
 274  
 275 **Fig. 3** The joint label of Kinetics-skeleton and NTU-RGB+D datasets.

276 The original paper of the NTU-RGB+D dataset recommended two benchmarks: 1) cross-  
 277 subject (X-Sub) benchmark: The dataset in this benchmark is divided into a training set (40,320  
 278 clips) and a validation set (16,560 clips). The subjects in these two subsets are different; 2) cross-  
 279 view (X-View) benchmark: The training set in this benchmark contains 37,920 clips captured by  
 280 cameras 2 and 3, and the validation set contains 18,960 clips captured by camera 1<sup>32</sup>. We follow  
 281 this convention and report the top-1 recognition accuracy of the two benchmarks.

282 Fig. 3(b) shows the joint label of the NTU-RGB+D dataset. The joint labels are: 1 base of the  
 283 spine, 2 middle of the spine, 3 neck, 4 head, 5 left shoulder, 6 left elbow, 7 left wrist, 8 left hand,  
 284 9 right shoulder, 10 right elbow, 11 right wrist, 12 right hand, 13 left hip, 14 left knee, 15 left

285 ankle, 16 left foot, 17 right hip, 18 right knee, 19 right ankle, 20 right foot, 21 spine, 22 left hand  
286 tip, 23 left hand Thumb, 24 right hand tip, 25 right thumb.

## 287 4.2 Effectiveness Analysis of AM-STGCN

288 In this section, we first conduct a lot of ablation experiments on the Kinetics-skeleton dataset: 1)  
289 Adding attention block after the ConvS (spatial convolution) of different layers of the ST-GCN;  
290 2) Adding multiple attention blocks after the ConvS of different layers; 3) Adding attention  
291 blocks after ConvT (temporal convolution) of the layer; 4) Adding two other attention  
292 mechanisms with different structures, CBAM<sup>29</sup> and SENet<sup>30</sup>, to ST-GCN. Experiments are then  
293 performed on NTU-RGB+D dataset to verify the generalization capabilities of the proposed  
294 model AM-STGCN.

### 295 4.2.1 Baseline

296 In order to evaluate the recognition performance of our improved model, we used baseline for  
297 comparison experiments. Since our model is improved on the basis of the ST-GCN model, we  
298 use the ST-GCN model as a baseline to analyze the advantages of AM-STGCN. We reproduced  
299 the ST-GCN model on the Kinetics dataset based on the Ref. 3, and obtained very close results to  
300 the original paper (see Table 1).

301

**Table 1** Baseline.

Method	Top-1(%)	Top-5(%)
ST-GCN <sup>3</sup>	30.7	52.8
Our ST-GCN Baseline	30.7	53.7

302 4.2.2 Ablation experiment

303 **Table 2** The results of adding one attention block to the different layers of the ST-GCN. ST-GCN1's ConvS + 1  
 304 represents adding one attention block after the ConvS (spatial convolution) of the first layer of the ST-GCN.  
 305 Thereafter, Tables 3, 4, 5, and 6 have the same representation rules.

Method	Top-1(%)	Top-5(%)
Our ST-GCN Baseline	30.7	53.7
ST-GCN <sub>1</sub> 's ConvS + 1	31.6	54.3
ST-GCN <sub>2</sub> 's ConvS + 1	<b>31.9</b>	<b>54.7</b>
ST-GCN <sub>3</sub> 's ConvS + 1	<b>31.9</b>	<b>54.7</b>
ST-GCN <sub>4</sub> 's ConvS + 1	31.3	53.8
ST-GCN <sub>9</sub> 's ConvS + 1	31.0	53.7

306 Table 2 shows the experimental results of adding one attention block after the ConvS (spatial  
 307 convolution) of different layers of the ST-GCN model. The results demonstrate that no matter  
 308 which layer we add an attention block to, the recognition accuracy always higher than the  
 309 baseline. The improvement of adding one attention block in the second and third layers is similar,  
 310 which can lead to  $\sim 1.2\%$  (on Top1) improvement over the baseline. The results of the remaining  
 311 layers are slightly lower.

313 **Table 3** The results of adding multiple attention blocks to different layers.

Method	Top-1(%)	Top-5(%)
Our ST-GCN Baseline	30.7	53.7
ST-GCN <sub>1</sub> 's ConvS + 2	32.0	54.5
ST-GCN <sub>2</sub> 's ConvS + 2	32.1	54.4
ST-GCN <sub>3</sub> 's ConvS + 2	31.4	54.4
ST-GCN <sub>1</sub> 's ConvS + 3	30.6	53.1

ST-GCN <sub>2</sub> 's ConvS + 3	31.1	53.5
ST-GCN <sub>3</sub> 's ConvS + 3	<b>32.2</b>	<b>55.1</b>
ST-GCN <sub>4</sub> 's ConvS + 3	31.1	53.1

314  
315 Table 3 shows the results of adding multiple attention blocks to different layers of the ST-  
316 GCN. It can be seen from Table 2 that adding one attention block to the first few layers of the  
317 model is better than adding to the lower layer, so in the experiment of Table 3, we add two and  
318 three attention blocks after the ConvS (spatial convolution) of the first few layers of ST-GCN.  
319 Obviously, the results of adding multiple attention blocks after ConvS of a layer outperform  
320 adding a single attention block, especially on ST-GCN<sub>3</sub>'s ConvS + 3, which can lead to 1.5%  
321 (on Top1) and 1.4% (on Top5) improvement over the baseline. It demonstrates that more  
322 attention blocks usually lead to better performance. We argue that multiple attention blocks can  
323 reinforce the correlation information learned in the previous attention block, thus assigning each  
324 node a more appropriate weight.

325 **Table 4** The results of adding multiple attention blocks to multi-layers.

Method	Top-1(%)	Top-5(%)
Our ST-GCN Baseline	30.7	53.7
ST-GCN <sub>2</sub> 's ConvS + 1 ST-GCN <sub>3</sub> 's ConvS + 1	31.4	54.1
ST-GCN <sub>1</sub> 's ConvS + 2 ST-GCN <sub>2</sub> 's ConvS + 2	30.9	53.3
ST-GCN <sub>2</sub> 's ConvS + 2 ST-GCN <sub>3</sub> 's ConvS + 2	<b>32.3</b>	<b>55.1</b>
ST-GCN <sub>1</sub> 's ConvS + 2 ST-GCN <sub>3</sub> 's ConvS + 2	31.5	54.2

326



327 Table 4 shows the results of adding multiple attention blocks to multi-layers of the ST-GCN  
 328 model. As shown in Tables 2, 3 and 4, we can find that only the third combination (ST-GCN<sub>2</sub>'s  
 329 ConvS + 2 & ST-GCN<sub>3</sub>'s ConvS + 2) improves accuracy compared to adding attention blocks to  
 330 single layer. The rest of the combinations do not improve accuracy compared to the individual  
 331 structure in the combination.

332 **Table 5** The results of adding attention blocks after ConvT (temporal convolution) of one layer.

Method	Top-1(%)	Top-5(%)
Our ST-GCN Baseline	30.7	53.7
ST-GCN <sub>2</sub> 's ConvT + 2	32.0	54.9
ST-GCN <sub>3</sub> 's ConvT + 3	<b>32.9</b>	<b>55.4</b>
ST-GCN <sub>5</sub> 's ConvT + 3	31.7	54.3

333 Table 5 shows the results of adding attention blocks after ConvT (temporal convolution) of  
 334 different layers of the ST-GCN model. Comparing the results of Table 3 and Table 5, we can  
 335 find that adding attention blocks after ConvT perform better than after ConvS. ST-GCN<sub>3</sub>'s  
 336 ConvT + 3 obtain the best improvement of adding attention blocks after ConvT, which  
 337 outperforms Our ST-GCN Baseline by 2.2% and 1.7% on Top1 and Top5 recognition accuracies;  
 338 ST-GCN<sub>3</sub>'s ConvS + 3 obtain the best improvement of adding attention blocks after ConvS,  
 339 which outperforms Our ST-GCN Baseline by 1.5% and 1.4% on Top1 and Top5 recognition  
 340 accuracies. One possible explanation is that ConvT has a bigger kernel size ( $9 \times 1$ ) and ConvS  
 341 has a small kernel size ( $1 \times 1$ ), thus ConvS is insufficient to capture precise spatial information.  
 342 Adding attention blocks after ConvT can learn the correlation of all nodes in all frames, while  
 343 adding attention blocks after ConvS can only learn the correlation of all nodes in one frame, thus  
 344 adding attention blocks after ConvT perform better than after ConvS.

346

**Table 6** The results of adding attention blocks after ConvT and ConvS of multi-layers.

Method	Top-1(%)	Top-5(%)
Our ST-GCN Baseline	30.7	53.7
ST-GCN <sub>2</sub> 's ConvT + 2 ST-GCN <sub>3</sub> 's ConvT + 3	<b>32.3</b>	<b>54.4</b>
ST-GCN <sub>2</sub> 's ConvS + 1 ST-GCN <sub>2</sub> 's ConvT + 2	31.5	53.8
ST-GCN <sub>2</sub> 's ConvS + 2 ST-GCN <sub>3</sub> 's ConvT + 3	31.8	54.0

347

348

Table 6 shows the results of adding attention blocks after ConvT and ConvS of multi-layers.

349

As shown in Tables 2, 3, 5 and 6, we can see that none of the combinations in Table 6 improves

350

accuracy compared to adding attention blocks to single layer. The results of Table 4 and 6 prove

351

that adding attention blocks to multiple layers does not further improve accuracy.

352

From Tables 2, 3, 4, 5 and 6, we find that adding attention blocks to the second and third

353

layer of ST-GCN can result in better performance. The possible reason is that the features

354

learned in these two layers are more consistent with the semantic representation of human

355

motion.

356

**Table 7** The results of adding CBAM and SENet to ST-GCN.

Method	Top-1(%)	Top-5(%)
ST-GCN+CBAM	31.9	54.3
ST-GCN+SENet	31.6	54.2
<b>Our AM-STGCN</b>	<b>32.9</b>	<b>55.4</b>

357

358

We selected two other attention mechanisms with different structures, CBAM<sup>29</sup> and SENet<sup>30</sup>,

359

to be added to ST-GCN. CBAM contains spatial attention and channel attention, while SENet is

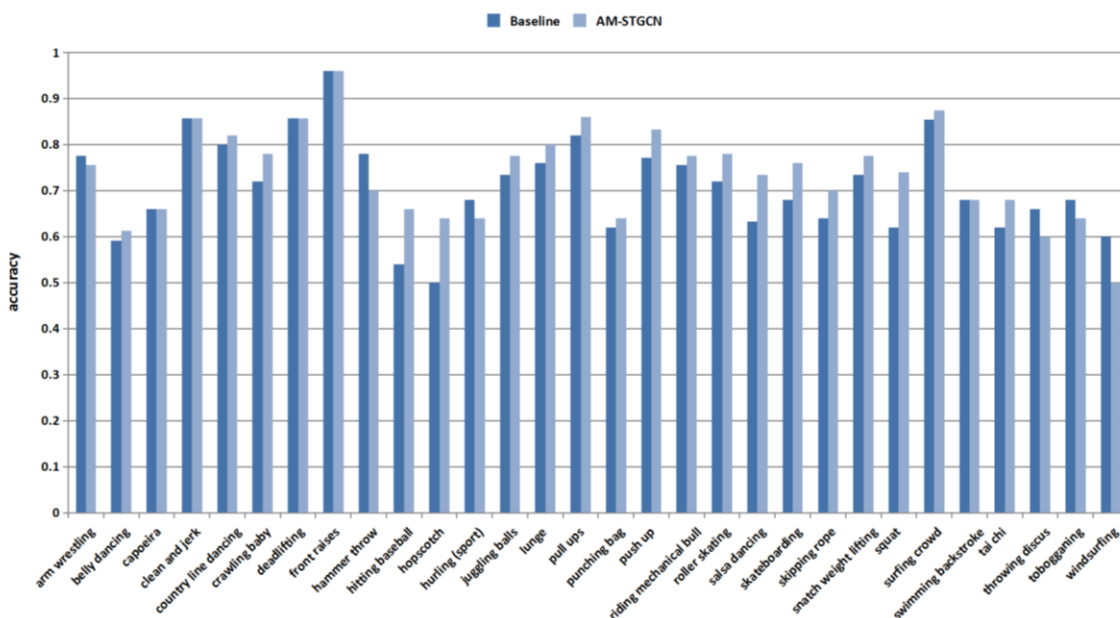
360

just channel attention. Table 7 shows the results of adding CBAM and SENet. As shown in

361 Table7, the results of our method are clearly better than those of the other two attention  
 362 structures, which prove that our attention mechanism is more suitable for ST-GCN.

363 4.2.3 Further analysis on “Kinetics-Motion”

364 The authors of ST-GCN select a subset of 30 classes strongly related with body motions, named  
 365 as “Kinetics-Motion<sup>3</sup>”. For a detailed comparison, we further investigate the per-class differences  
 366 in accuracy on this subset. In Fig. 4, the horizontal axis is the action category of “Kinetics-  
 367 Motion”, and the vertical axis is the accuracy of per-class. The dark blue represents Our ST-  
 368 GCN Baseline and the light blue represents AM-STGCN, here AM-STGCN is the optimal  
 369 structure (i.e., ST-GCN<sub>3</sub>’s ConvT + 3) obtained after the analysis in the previous section. It can  
 370 be observed obviously that the accuracy of most actions get improved. Some classes even get  
 371 more than 10% improvement, such as hitting baseball, hopscotch, salsa dancing and squat. These  
 372 results also verify the superiority of our model for skeleton-based action recognition, in  
 373 particular on those classes strongly related with body motions.



374

375

**Fig. 4** Category accuracies on the “Kinetics Motion” subset of the Kinetics dataset.

376 4.2.4 Time comparison on Kinetics

377 The Kinetics dataset provides a training set of 240,000 video clips, each clip contain 300 frames.  
378 Every frame of the video clips is converted into a sequence of human skeletons represented by  
379 coordinates through OpenPose<sup>34</sup> toolbox. We compared the training time of one epoch of AM-  
380 STGCN model and our ST-GCN baseline on Kinetics dataset, and the results are shown in Table  
381 8. ST-GCN<sub>3</sub>'s ConvS + 3 and ST-GCN<sub>3</sub>'s ConvT + 3, which performed better in the above  
382 experiments, are selected to be compared with our ST-GCN baseline. The training time of ST-  
383 GCN<sub>3</sub>'s ConvS + 3 and our ST-GCN baseline are similar, and ST-GCN<sub>3</sub>'s ConvT adds the  
384 calculation in temporal dimension, so the training time is a little longer. These results  
385 demonstrate that our AM-STGCN model do not add much time cost than ST-GCN model.

386 **Table 8** The training time of AM-STGCN and ST-GCN methods.

Method	The number of skeleton sequence.	Training time of one epoch. (h)
Our ST-GCN Baseline	240,000	0.58
ST-GCN <sub>3</sub> 's ConvS + 3	240,000	0.61
ST-GCN <sub>3</sub> 's ConvT + 3	240,000	0.70

387

388 4.2.5 Comparison with state-of-the-art methods

389 On Kinetics dataset, we compare AM-STGCN with “Feature Encoding”<sup>10</sup>, Deep LSTM<sup>32</sup>,  
390 Temporal ConvNet<sup>14</sup> and ST-GCN<sup>3</sup> methods. Their recognition performance in terms of Top-1  
391 and Top-5 accuracies are listed in Table 9. Obviously, our AM-STGCN with using attention  
392 module outperforms ST-GCN by 2.2% and 2.6% on Top1 and Top5 recognition accuracies  
393 respectively. It can be seen from Table 9 that our AM-STGCN is able to outperform previous  
394 representative methods.

395

**Table 9** Comparison with the state-of-the-art on Kinetics dataset.

Method	Date	Top-1(%)	Top-5(%)
Feature Encoding <sup>10</sup>	2015	14.9	25.8
Deep LSTM <sup>32</sup>	2016	16.4	35.3
Temporal ConvNet <sup>14</sup>	2017	20.3	40.0
ST-GCN <sup>3</sup>	2018	30.7	52.8
Our ST-GCN Baseline	-	30.7	53.7
<b>Our AM-STGCN</b>	-	<b>32.9</b>	<b>55.4</b>

396

397 We found that most of the current skeleton-based action recognition studies are conducted on  
 398 NTU-RGB+D dataset, so we compare our method with state-of-the-art methods on NTU-  
 399 RGB+D dataset.

400 On NTU-RGB+D dataset, we compare AM-STGCN with Lie Group<sup>9</sup>, H-RNN<sup>11</sup>, Deep  
 401 LSTM<sup>32</sup>, VA-LSTM<sup>13</sup>, Temporal ConvNet<sup>14</sup>, Two-stream CNN<sup>16</sup>, HCN<sup>17</sup>, STA-LSTM<sup>12</sup>, GCA-  
 402 LSTM<sup>18</sup>, ARRN-LSTM<sup>20</sup>, MANs<sup>19</sup>, ST-GCN<sup>3</sup>, DPRL+GCNN<sup>26</sup>, SR-TSL<sup>22</sup>, PB-GCN<sup>24</sup> and  
 403 AGCN<sup>23</sup> methods. The results are shown in Table 10.

404 **Comparisons with hand-craft feature based methods, CNN based methods and RNN  
 405 based methods.** Table 10 shows that the performance of graph convolution based methods is  
 406 generally better than hand-craft feature based methods, CNN based methods and RNN based  
 407 methods. In particular, our AM-STGCN obtains very close results to HCN method on cross-  
 408 view (X-View) benchmark, which performs best among CNN based methods. At the same time,  
 409 multi-person feature fusion is added in HCN, thus resulting in better performance on cross-  
 410 subject (X-Sub) benchmark, but it also leads to the increase of computation.

411 **Comparisons with other methods based on attention.** We compare AM-STGCN with  
 412 other methods based on attention including STA-LSTM<sup>12</sup>, GCA-LSTM<sup>18</sup>, ARRN-LSTM<sup>20</sup> and

413 MANs<sup>19</sup>. From Table 10, we can see that our AM-STGCN is better than any other result except  
414 for MANs under the X-View benchmark. MANs consists of Temporal Attention Recalibration  
415 Module (TARM) and DenseNet-161, we can find that their baseline is higher than ST-GCN,  
416 which may be due to DenseNet-161, because DenseNet-161 is much deeper and more complex  
417 than ST-GCN. On X-View benchmark, our AM-STGCN outperforms ST-GCN by 3.1% and  
418 MANs outperforms MANs (no attention) by 1.07%, which prove that our method can improve  
419 the performance of the model more.

420 **Comparisons with graph convolution based methods.** 1) Single stream network. In Table  
421 10, we can see clearly that our AM-STGCN with using attention module outperforms ST-GCN  
422 by 1.9% and 3.1% on cross-view (X-View) benchmark and cross-subject (X-Sub) benchmark  
423 respectively, which prove that our AM-STGCN model is equally effectiveness on NTU-RGB+D  
424 dataset. Our AM-STGCN performs very close results to DPRL+GCNN on cross-subject (X-Sub)  
425 benchmark and outperforms DPRL+GCNN by 1.6% on cross-view (X-View) benchmark in  
426 Table 10. 2) Two-stream networks. The joint locations is the only input data of our AM-STGCN.  
427 SR-TSL, PB-GCN and AGCN all have another form of input data as input to different streams,  
428 thus forming a two-stream networks. SR-TSL(Position), PB-GCN(*Jloc*) and Js-AGCN are the  
429 same as ST-GCN with only joint locations as input data. Among these methods, it can be seen  
430 obviously from Table 10 that our AM-STGCN is superior to SR-TSL(Position) and PB-  
431 GCN(*Jloc*) on both cross-subject (X-Sub) and cross-view (X-View) benchmark. In the paper of  
432 AGCN, we find AGCN's baseline is 92.7% on cross-view (X-View) benchmark, outperforms  
433 ST-GCN by 4.4%, but Js-AGCN outperforms their baseline by only 1%. We think it may be that  
434 different experimental environments cause different baselines. So in terms of relative increase in  
435 accuracy, our method has achieved a good performance improvement. In addition, we have

436 added our attention module to Js-AGCN. In Table 10, the results of Js-AGCN+our attention  
 437 outperforms Our Js-AGCN Baseline by 0.5% and 0.4% on cross-view (X-View) benchmark and  
 438 cross-subject (X-Sub) benchmark respectively, which shows that our attention mechanism is also  
 439 effective on AGCN method, and proves that our method has certain robustness.

440 These results show our AM-STGCN model achieves a significant performance improvement.

441 **Table 10** Comparison with the state-of-the-art on NTU-RGB+D dataset.

Method	Date	X-Sub(%)	X-View(%)
Lie Group <sup>9</sup>	2014	50.1	52.8
H-RNN <sup>11</sup>	2015	59.1	64.0
Deep LSTM <sup>32</sup>	2016	60.7	67.3
Temporal ConvNet <sup>14</sup>	2017	74.3	83.1
VA-LSTM <sup>13</sup>	2017	79.4	87.6
Two-stream CNN <sup>16</sup>	2017	83.2	89.3
HCN <sup>17</sup>	2018	86.5	91.1
STA-LSTM <sup>12</sup>	2017	73.4	81.2
GCA-LSTM <sup>18</sup>	2017	74.4	82.8
ARRN-LSTM <sup>20</sup>	2019.04	81.8	89.6
MANs (no attention) <sup>19</sup>	2018	81.41	92.15
MANs <sup>19</sup>		83.01	93.22
ST-GCN <sup>3</sup>	2018	81.5	88.3
DPRL+GCNN <sup>26</sup>	2018	83.5	89.8
SR-TSL(Position) <sup>22</sup>		78.8	88.2
SR-TSL(Velocity) <sup>22</sup>	2018	82.2	90.6
SR-TSL <sup>22</sup>		84.8	92.4

PB-GCN( $J_{loc}$ ) <sup>24</sup>		82.8	90.3
PB-GCN( $D_R  D_T$ ) <sup>24</sup>	2018	87.5	93.2
Js-AGCN <sup>23</sup>		-	93.7
Bs-AGCN <sup>23</sup>	2019.05	-	93.2
2s-AGCN <sup>23</sup>		88.5	95.1
Our Js-AGCN Baseline	-	85.9	93.7
Js-AGCN + our attention	-	86.4	94.1
<b>Our AM-STGCN</b>	-	<b>83.4</b>	<b>91.4</b>

442

## 443 5 Conclusion

444 In this paper, we propose a new skeleton-based action recognition method called attention  
445 module-based Spatial Temporal Graph Convolutional Networks(AM-STGCN), which can  
446 overcome the weakness of ST-GCN model. In order to capture global information of skeleton  
447 sequences, attention modules are added to learn the correlation information between all joints of  
448 both spatial and temporal dimension. So AM-STGCN can extract long-range relationships from  
449 input skeleton sequences, which improve the ability to model the dynamic change of human  
450 body motions. Experiments on two large-scale action recognition datasets Kinetics and NTU-  
451 RGB+D achieve the better results, which indicate that AM-STGCN can effectively improve the  
452 recognition accuracy. In future, we will improve our AM-STGCN in many possible directions,  
453 such as improving attention modules or merging RGB modality.

454



455 **Acknowledgments**

456 This work is supported by the National Natural Science Foundation of China (Grant Nos.  
457 61871182, 61302163), Hebei Province Natural Science Foundation (Grant Nos.F2015502062),  
458 Hebei province science and technology support (Grant Nos.13210905).

459

460 **References**

- 461 1. H. Wang and C. Schmid, "Action recognition with improved trajectories," in *2013 IEEE*  
462 *International Conference on Computer Vision*, pp. 3551-3558, IEEE, Sydney, NSW (2013)  
463 [[doi:10.1109/ICCV.2013.441](https://doi.org/10.1109/ICCV.2013.441)].
- 464 2. H. Wang et al., "Dense trajectories and motion boundary descriptors for action recognition," in  
465 *International journal of computer vision*, **103**(1), 60-79 (2013).
- 466 3. S. Yan, Y. Xiong, and D. Lin, "Spatial temporal graph convolutional networks for skeleton-based  
467 action recognition," in *Thirty-Second AAAI Conference on Artificial Intelligence*, pp. 7444-7452,  
468 AAAI Press, New Orleans, Louisiana, USA (2018).
- 469 4. K. Simonyan, A. Zisserman, "Two-stream convolutional networks for action recognition in videos,"  
470 in *Advances in neural information processing systems*, pp. 568-576 (2014).
- 471 5. C. Feichtenhofer, A. Pinz and A. Zisserman, "Convolutional two-Stream network fusion for video  
472 action recognition," in *Proceedings of the IEEE Conference on Computer Vision and Pattern*  
473 *Recognition*, pp. 1933-1941, IEEE, Las Vegas, NV (2016) [[doi:10.1109/CVPR.2016.213](https://doi.org/10.1109/CVPR.2016.213)].
- 474 6. D. Tran et al., "Learning spatiotemporal features with 3D convolutional networks," in *2015 IEEE*  
475 *International Conference on Computer Vision*, pp. 4489-4497, IEEE, Santiago (2015)  
476 [[doi:10.1109/ICCV.2015.510](https://doi.org/10.1109/ICCV.2015.510)].

- 477 7. J. Carreira and A. Zisserman, “Quo Vadis, Action Recognition? A New Model and the Kinetics  
478 Dataset,” in *2017 IEEE Conference on Computer Vision and Pattern Recognition*, pp. 4724-4733,  
479 IEEE, Honolulu, HI (2017) [[doi:10.1109/CVPR.2017.502](https://doi.org/10.1109/CVPR.2017.502)].
- 480 8. W. Du, Y. Wang and Y. Qiao, “RPAN: An end-to-end recurrent pose-attention network for action  
481 recognition in videos,” in *2017 IEEE International Conference on Computer Vision*, pp. 3745-3754,  
482 IEEE, Venice (2017) [[doi:10.1109/ICCV.2017.402](https://doi.org/10.1109/ICCV.2017.402)].
- 483 9. R. Vemulapalli, F. Arrate, and R. Chellappa, “Human action recognition by representing 3d skeletons  
484 as points in a lie group,” in *Proceedings of the IEEE Conference on Computer Vision and Pattern  
485 Recognition*, pp. 588-595, IEEE, Columbus, OH (2014) [[doi:10.1109/CVPR.2014.82](https://doi.org/10.1109/CVPR.2014.82)].
- 486 10. B. Fernando et al., “Modeling video evolution for action recognition,” in *Proceedings of the IEEE  
487 Conference on Computer Vision and Pattern Recognition*, pp. 5378-5387, IEEE, Boston, MA (2015)  
488 [[doi:10.1109/CVPR.2015.7299176](https://doi.org/10.1109/CVPR.2015.7299176)].
- 489 11. Y. Du, W. Wang, and L. Wang, “Hierarchical recurrent neural network for skeleton based action  
490 recognition,” in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*,  
491 pp. 1110-1118, IEEE, Boston, MA (2015) [[doi:10.1109/CVPR.2015.7298714](https://doi.org/10.1109/CVPR.2015.7298714)].
- 492 12. S. Song et al., “An end-to-end spatio-temporal attention model for human action recognition from  
493 skeleton data,” in *Proceedings of the Thirty-First AAAI Conference on Artificial Intelligence*, pp.  
494 4263-4270, AAAI Press, San Francisco, California, USA (2017).
- 495 13. P. Zhang et al., “View adaptive recurrent neural networks for high performance human action  
496 recognition from skeleton data,” in *2017 IEEE International Conference on Computer Vision*, pp.  
497 2136-2145, IEEE, Venice (2017) [[doi:10.1109/ICCV.2017.233](https://doi.org/10.1109/ICCV.2017.233)].
- 498 14. T. S. Kim and A. Reiter, “Interpretable 3d human action analysis with temporal convolutional  
499 networks,” in *IEEE Conference on Computer Vision and Pattern Recognition Workshops (CVPRW)*,  
500 pp. 1623-1631, IEEE, Honolulu, HI (2017) [[doi:10.1109/CVPRW.2017.207](https://doi.org/10.1109/CVPRW.2017.207)].
- 501 15. M. Liu, H. Liu, and C. Chen, “Enhanced skeleton visualization for view invariant human action  
502 recognition,” in *Pattern Recognition*, **68**, pp. 346-362, Elsevier (2017).

- 503 16. C. Li et al., “Skeleton-based action recognition with convolutional neural networks,” in *2017 IEEE*  
504 *International Conference on Multimedia & Expo Workshops (ICMEW)*, pp. 597-600, IEEE, Hong  
505 Kong (2017) [[doi:10.1109/ICMEW.2017.8026285](https://doi.org/10.1109/ICMEW.2017.8026285)].
- 506 17. C. Li et al., “Co-occurrence feature learning from skeleton data for action recognition and detection  
507 with hierarchical aggregation,” in *Proceedings of the 27th International Joint Conference on*  
508 *Artificial Intelligence*, pp. 786-792, AAAI Press, Stockholm, Sweden (2018).
- 509 18. J. Liu et al., “Global Context-Aware Attention LSTM Networks for 3D Action Recognition,” in *2017*  
510 *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 3671-3680, IEEE,  
511 Honolulu, HI (2017) [[doi: 10.1109/CVPR.2017.391](https://doi.org/10.1109/CVPR.2017.391)].
- 512 19. C. Xie et al., “Memory attention networks for skeleton-based action recognition,” in *Proceedings of*  
513 *the 27th International Joint Conference on Artificial Intelligence*, pp. 1639-1645, AAAI Press,  
514 Stockholm, Sweden (2018).
- 515 20. W. Zheng et al., “Relational Network for Skeleton-Based Action Recognition,” arXiv preprint  
516 [arXiv:1805.02556v4](https://arxiv.org/abs/1805.02556v4), 2019.
- 517 21. C. Li et al., “Spatio-temporal graph convolution for skeleton based action recognition,” in *Thirty-*  
518 *Second AAAI Conference on Artificial Intelligence*, pp. 3482-3489, AAAI Press, New Orleans,  
519 Louisiana, USA (2018).
- 520 22. C. Si et al., “Skeleton-based action recognition with spatial reasoning and temporal stack learning,” in  
521 *Proceedings of the European Conference on Computer Vision (ECCV)*, Lecture Notes in Computer  
522 Science, vol **11205**, pp. 106-121, Springer, Cham (2018) [[https://doi.org/10.1007/978-3-030-01246-](https://doi.org/10.1007/978-3-030-01246-5_7)  
523 [5\\_7](https://doi.org/10.1007/978-3-030-01246-5_7)].
- 524 23. L. Shi et al., “Two-Stream Adaptive Graph Convolutional Networks for Skeleton-Based Action  
525 Recognition,” in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*,  
526 pp.12026-12035 (2019).
- 527 24. K. Thakkar, P.J. Narayanan, “Part-based Graph Convolutional Network for Action Recognition,”  
528 arXiv preprint [arXiv:1809.04983](https://arxiv.org/abs/1809.04983), 2018.

- 529 25. X. Gao et al., “Optimized Skeleton-based Action Recognition via Sparsified Graph Regression,”  
530 arXiv preprint arXiv:1811.12013v2, 2019.
- 531 26. Y. Tang et al., “Deep Progressive Reinforcement Learning for Skeleton-Based Action Recognition,”  
532 in *2018 IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 5323-5332, IEEE,  
533 Salt Lake City, UT (2018) [[doi: 10.1109/CVPR.2018.00558](https://doi.org/10.1109/CVPR.2018.00558)].
- 534 27. X. Wang et al., “Non-local neural networks,” in *Proceedings of the IEEE Conference on Computer*  
535 *Vision and Pattern Recognition*, pp. 7794-7803, IEEE, Salt Lake City, UT, USA (2018)  
536 [[doi:10.1109/CVPR.2018.00813](https://doi.org/10.1109/CVPR.2018.00813)].
- 537 28. Y. Du et al., “Interaction-Aware Spatio-Temporal Pyramid Attention Networks for Action  
538 Classification,” in *Proceedings of the European Conference on Computer Vision (ECCV)*, Lecture  
539 Notes in Computer Science, vol **11220**, pp. 388-404, Springer, Cham (2018)  
540 [[https://doi.org/10.1007/978-3-030-01270-0\\_23](https://doi.org/10.1007/978-3-030-01270-0_23)].
- 541 29. S. Woo et al., “CBAM: Convolutional Block Attention Module,” in *Proceedings of the European*  
542 *Conference on Computer Vision (ECCV)*, Lecture Notes in Computer Science, vol **11211**, pp. 3-19,  
543 Springer, Cham (2018) [[https://doi.org/10.1007/978-3-030-01234-2\\_1](https://doi.org/10.1007/978-3-030-01234-2_1)].
- 544 30. J. Hu, L. Shen, and G. Sun, “Squeeze-and-Excitation Networks,” in *2018 IEEE/CVF Conference on*  
545 *Computer Vision and Pattern Recognition*, pp. 7132-7141, IEEE, Salt Lake City, UT (2018)  
546 [[doi:10.1109/CVPR.2018.00745](https://doi.org/10.1109/CVPR.2018.00745)].
- 547 31. W. Kay et al., “The kinetics human action video dataset,” arXiv preprint arXiv:1705.06950, 2017.
- 548 32. A. Shahroudy et al., “NTU RGB+D: A large scale dataset for 3D human activity analysis,”  
549 in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pp. 1010-1019,  
550 IEEE, Las Vegas, NV (2016) [[doi:10.1109/CVPR.2016.115](https://doi.org/10.1109/CVPR.2016.115)].
- 551 33. K. He et al., “Deep residual learning for image recognition,” in *Proceedings of the IEEE Conference*  
552 *on Computer Vision and Pattern Recognition*, pp. 770-778, IEEE, Las Vegas, NV (2016)  
553 [[doi:10.1109/CVPR.2016.90](https://doi.org/10.1109/CVPR.2016.90)].

554 34. Z. Cao et al., “Realtime multi-person 2D pose estimation using part affinity fields,” in *Proceedings*  
555 *of the IEEE Conference on Computer Vision and Pattern Recognition*, pp. 1302-1310, IEEE,  
556 Honolulu, HI (2017) [[doi:10.1109/CVPR.2017.143](https://doi.org/10.1109/CVPR.2017.143)].

557

558 **Yinghui Kong** is a professor at North China Electric Power University, Baoding, China. She  
559 received BS degree in Communication Engineering from Xidian University, Xian, China, in  
560 1987, and MS degree in Power System & Its Automation from North China Electric Power  
561 University, Beijing, China, in 1993, and PhD degree in Electrical Theory and New Techniques  
562 from North China Electric Power University in 2009. Her research interests include deep  
563 learning, computer vision, behavior recognition, expression recognition.

564

565

## 566 **Caption List**

567

568 **Fig. 1** (a) Spatial-temporal graph of the skeleton and (b) Partitioning strategy, different colors  
569 represent different subsets.

570 **Fig. 2** The structure of AM-STGCN.

571 **Fig. 3** The joint label of Kinetics-skeleton and NTU-RGB+D datasets.

572 **Fig. 4** Category accuracies on the “Kinetics Motion” subset of the Kinetics dataset.

573 **Table 1** Baseline.

574 **Table 2** The results of adding a attention block to the different layers of the ST-GCN. ST-  
575 GCN1’s ConvS + 1 represents adding one attention block after the ConvS of the first layer of the  
576 ST-GCN. Thereafter, Tables 3, 4, 5, and 6 have the same representation rules.

577 **Table 3** The results of adding multiple attention blocks to different layers.

578 **Table 4** The results of adding multiple attention blocks to multi-layer.

- 579 **Table 5** The results of adding attention blocks after ConvT of one layer.
- 580 **Table 6** The results of adding attention blocks after ConvT and ConvS of multi-layer.
- 581 **Table 7** The results of adding CBAM and SENet to ST-GCN.
- 582 **Table 8** The training time of AM-STGCN and STGCN methods.
- 583 **Table 9** Comparison with the state-of-the-art on Kinetics dataset.
- 584 **Table 10** Comparison with the state-of-the-art on NTU-RGB+D dataset.

