

Action-Attending Graphic Neural Network

Chaolong Li*, Zhen Cui*, *Member, IEEE*, Wenming Zheng, *Member, IEEE*, Chunyan Xu, *Member, IEEE*, Rongrong Ji, *Senior Member, IEEE*, and Jian Yang, *Member, IEEE*

Abstract—The motion analysis of human skeletons is crucial for human action recognition, which is one of the most active topics in computer vision. In this paper, we propose a fully end-to-end action-attending graphic neural network (A²GNN) for skeleton-based action recognition, in which each irregular skeleton is structured as an undirected attribute graph. To extract high-level semantic representation from skeletons, we perform the local spectral graph filtering on the constructed attribute graphs like the standard image convolution operation. Considering not all joints are informative for action analysis, we design an action-attending layer to detect those salient action units (AUs) by adaptively weighting skeletal joints. Herein the filtering responses are parameterized into a weighting function irrelevant to the order of input nodes. To further encode continuous motion variations, the deep features learnt from skeletal graphs are gathered along consecutive temporal slices and then fed into a recurrent gated network. Finally, the spectral graph filtering, action-attending and recurrent temporal encoding are integrated together to jointly train for the sake of robust action recognition as well as the intelligibility of human actions. To evaluate our A²GNN, we conduct extensive experiments on four benchmark skeleton-based action datasets, including the large-scale challenging NTU RGB+D dataset. The experimental results demonstrate that our network achieves the state-of-the-art performances.

Index Terms—Human action recognition, Skeleton-based action recognition, Convolutional neural networks, Attention mechanism.

I. INTRODUCTION

HUMAN action recognition is an active area of research with wide applications such as video surveillance, games console, robot vision, etc. Over the past few decades, human action recognition from 2D RGB video sequences has been extensively studied [1]. However, 2D cameras cannot fully capture human motions, which are actually located in 3D space. With the advent of depth sensors such as Microsoft Kinect and Asus Xtion PRO LIVE, 3D action recognition arises more attention of researchers due to its several advantages in segmenting foreground/background, resisting illumination variations, etc.

Recently several approaches have been developed to deal with 3D human action recognition. Generally they fall into two categories: depth map based or skeleton based. The depth map based methods directly extract volumetric and temporal features of overall point set from depth map sequences. The

skeleton based methods utilize 3D coordinates of skeletal joints estimated from depth maps to model actions of human body. As human body can be viewed as an articulated system of rigid bones connected by hinged joints, the actions of human body principally reflect in human skeletal motions in the 3D space [2]. Thereby, to model human skeletal joints should be more effective for action recognition, as suggested in the early work [3].

Skeleton based methods [4], [5], [6], [7], [8], [9], [10], [11], [12], [13] have become prevalent for human action recognition since Shotton *et al.* [14] successively estimated 3D coordinates of skeletal joints from depth maps. To represent trajectories of human body motions, the position, speed or even acceleration of skeletal joints are often gathered from several consecutive frames, and then modeled with some statistic methods (*e.g.*, histogram). Further, the skeletal joints may be partitioned into several parts according to the concurrent function of adjacent joints, so human actions are represented with the motion parameters of body parts. To encode temporal motions of skeletal joints, linear/non-linear dynamical systems are often used to model human action, *e.g.*, Hidden Markov Models (HMMs) [12] and Long Short Term Memory (LSTM) [15]. However, there still exist some key issues need to be studied deeply. First, how to extract robust high-level semantic features from irregular skeletons? The current skeleton-based methods usually employ some simple statistical features on skeletal joints or body parts. It is insufficient to describe and abstract spatial structure of skeleton at each time slice, whereas non-linear dynamic systems are only used to abstract high-level motion information. Second, which/what action units (AUs) identify a special human action? Most motions are produced from only a few joints or body parts, *e.g.*, waving with arm and hand, kicking ball with leg and foot. Detecting salient action units should be helpful to eliminate some useless motion noises as well as provide a cognitive explanation to understand human action.

To address the above problems, in this paper, we propose a fully end-to-end action-attending graphic neural network (A²GNN) for skeleton-based action recognition. To extract high-level semantic information from spatial skeletons, we represent each human skeleton with an undirected attribute graph, and then perform the local spectral filtering on the structured graph to laywisely abstract spatial skeletal information like the classic convolutional neural network (CNN) [16]. Hence, different from the recent graph-based method [13], which only took the traditional technique line of graph matching, A²GNN should be a true deep network directly learning from irregular skeletal structure. To detect those salient action units, we design an action-attending layer to adaptively weight skeletal joints for different human motions. Specifically, we

* C. Li and Z. Cui have equal contributions.

C. Li and W. Zheng are with the Key Laboratory of Child Development and Learning Science (Southeast University), Ministry of Education, School of Biological Science & Medical Engineering, Southeast University, Nanjing 210096, China (e-mail: {lichaolong, wenming_zheng}@seu.edu.cn).

Z. Cui, C. Xu and J. Yang are with the School of Computer Science and Engineering, Nanjing University of Science and Technology, Nanjing 210094, China (e-mail: {zhen.cui, cyx, csjyang}@njust.edu.cn).

R. Ji is with the School of Information Science and Engineering, Xiamen University, Xiamen 361005, China (e-mail: rrji@xmu.edu.cn).

parameterize the weighting process as one dynamic function, which takes the responses of local spectral filtering on skeletal graphs as the input. Incidentally, we draw a conclusion that the weighting process is irrelevant to the input order of skeletal joints, which thus can be used for those general graph tasks. After extracting high-level discriminant features at each time slice, we finally stack a recurrent neural network on a consecutive sequence to model temporal variations of different human actions. All processes including spectral graph filtering, action unit detection, temporal motion modeling are integrated into one network framework to train jointly. To evaluate our proposed method, we conduct extensive experiments on four benchmark skeleton-based action datasets: the Motion Capture Database HDM05 [17], the Florence 3D Action dataset [18], the Large Scale Combined (LSC) dataset [19], the NTU RGB+D dataset [20]. The experimental results demonstrate our A²GNN is competitive over the state-of-the-art methods.

In summary, our contributions are three folds:

- 1) Propose a fully end-to-end graphical neural network framework to deal with skeleton-based action recognition, where we model human skeletons as attribute graphs and then introduce spectral graph filtering to extract high-level skeletal features, .
- 2) Design a weighting way to adaptively detect salient action units for different human actions, which not only promotes human action recognition accuracy but also favors our cognitive understanding to human actions.
- 3) Achieve the state-of-the-art performances on the four benchmark datasets including the two large-scale challenging datasets: LSC and NTU RGB+D.

The reminder of this paper is organized as follows. In Section II, we introduce related work on skeleton based action recognition. In Section III, we first give an overview of our proposed network and introduce three main modules respectively. Implementation details of our proposed network are stated in Section IV. Experimental results and discussion are presented in Section V. Finally, we conclude this work in Section VI.

II. RELATED WORK

The most related works to ours are those methods of skeleton-based human action recognition. Here we briefly review them from the view of skeletal representation, including low-level statistical features and high-level semantic features.

Various low-level statistical features have been used to describe skeleton data in the past years. Generally they contain three categories: joint based, body part based and pose based features. Hussein *et al.* [21] used covariance matrix of skeletal joint coordinates over time to describe motion trajectories. In the literature [22], histogram of oriented displacement was used to represent the 3D trajectories of body joints. Ofli *et al.* [23] chose a few most informative joints, and employed some physical interpretable measures, such as the mean or variance of joint angles, maximum angular velocity of joints and so on, to encode skeletal motions. Considering joint-associated movements, Chaudhry *et al.* [24] divided human skeleton into several smaller parts hierarchically and depicted

each part with certain bio-inspired shapes. Vemulapalli *et al.* [10] formulated each body part over times into a curved manifold and performed action recognition in the Lie group. In view of skeletal postures, Zanfir *et al.* [25] proposed the moving pose descriptor by considering position, speed and acceleration information of body joints. Xia *et al.* [12] represented postures as histograms of 3D skeletal joint locations by casting selected joints into corresponding spatial histogram bins. These low-level statistical features are usually limited in representing those complex human actions.

High-level semantic features are popular for encoding human actions due to the robust representation ability. Especially, the temporal pyramid is used to hierarchically encode temporal dynamical variations, and all level features are gathered together to model actions [9], [10], [22], [26], [27]. To model the temporal dynamics, Vemulapalli *et al.* [10] used dynamic time warping (DTW) and Fourier temporal pyramid (FTP); Wang *et al.* [27] employed FTP to encoded local occupancy patterns over times; Chaudhry *et al.* [24] employed linear dynamical systems (LDSs) to learn the dynamical variation of features; Xia *et al.* [12] used HMMs to capture the temporal action dynamics. In addition, more recently, recurrent neural network has been employed to encode the temporal dynamic variations [15], [28]. However, these methods mainly focus on the deep encoding of temporal motions rather than spatial layout of joints.

Until recently, graph was used to represent skeletal motion in the literature [13], although a few graph-based algorithms [29], [30] had been proposed for action recognition on RGB videos. These graph-based methods usually took the traditional graph matching strategy after modeling skeletal joints or body parts into the graph structure. Therefore, two crucial issues, the construction of graph and the definition of graph kernel, need be conducted in these graph-based methods. Different from them, we perform spectral graph filtering on skeleton-induced graphs to extract high-level skeletal features from the spatial graphs. With a combination of recurrent motion encoding, the spatial-temporal features of skeletal motions are abstracted from the constructed end-to-end network. In contrast to the existing skeleton-based action recognition methods, another one important difference is that an adaptive action-attending mechanism is introduced to detect those salient action units w.r.t different human actions, which can benefit the final human action recognition.

III. THE PROPOSED A²GNN

In this section, we first give an overview of our proposed A²GNN, then we respectively introduce three involved modules: the learning of deep graphical features, the detection of salient action units and the dynamic modeling of temporal motions.

A. Overview

An overview of our proposed A²GNN is illustrated in Fig. 1. The input is a motion sequence of skeletons, in which each skeletal joint is described as a 3D coordinate (x, y, z) . For each skeleton at one time slice, we model it into an undirected

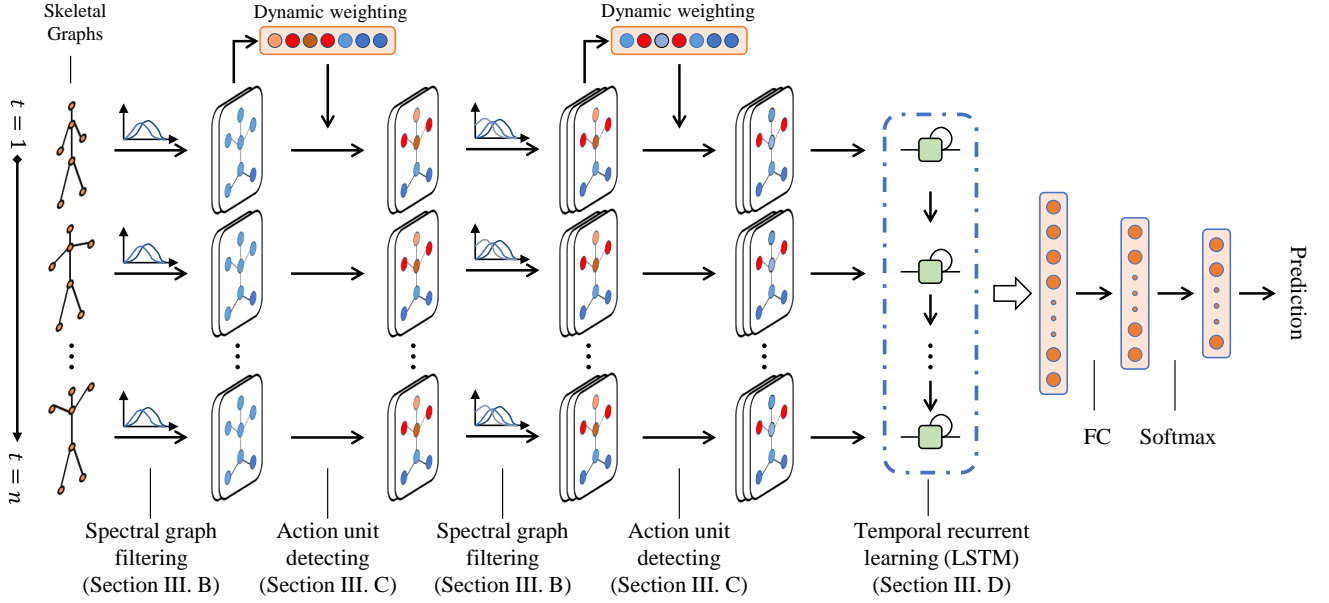


Fig. 1. The illustration of our proposed A²GNN architecture. An overall introduction can be found in Section III-A.

attribute graph, where each skeletal joint is one node and the bone between two joints is considered as one connected edge. For the weight between two nodes, we assign a connected edge to 1, otherwise 0. In addition, several alternative ways to weight edges are permissible, *e.g.*, Gaussian kernel. Another important characteristic is that each node is associated with an observed signal vector (a.k.a. attribute), which is 3D spatial coordinates of the joint. To reduce individual differences of different samples, we normalize the input signals with coordinating, scaling and rotating transformations.

For the constructed spatial skeletal graphs, in order to extract high-level skeletal features, we expect to perform convolutional filtering on them like on regular grid-shape images. To the end, we introduce local spectral filtering on graphs, inspired by signal theory on graphs [31] and the recent graph convolution [32]. To avoid the eigenvalue decomposition of Laplacian matrix, the original solution is approximated by a polynomial of Laplacian matrix, in which each k -order term derives a k -hop neighborhood subgraph like a local receptive field. The details are introduced in Section III-B.

Specially, we design an action-attending layer to detect salient action units by adaptively weighting skeletal joints for different human actions. Usually, specific action units are activated for different human actions, *e.g.*, hands and arms segments for clapping, hand, arm and head segments for drinking. Hence, weighting skeletal joints may reduce the disturbance of those useless joints, and benefit the final action recognition. The related details may be found in Section III-C.

By alternately stacking the spectral graph filtering layer and the action-attending layer, we may obtain the graph features of skeletal joints. After concatenating features of all joints at one time slice, we fed it into a recurrent network (LSTM) to encode feature variations at all temporal slices. Please refer to more details in Section III-D. Finally, a fully connected layer is used to gather the outputs of recurrent network and learn

the skeleton sequence representation followed by a softmax layer for classification. All processes including spectral graph filtering, action unit detection, temporal motion modeling are integrated into a network framework and jointly trained.

B. Learning Deep Graphical Features

We model a body skeleton into an undirected attribute graph $\mathcal{G} = (\mathcal{V}, \mathbf{A}, \mathbf{X})$ of N nodes, where $\mathcal{V} = \{v_1, \dots, v_N\}$ is the set of skeletal joints, \mathbf{A} is the (weighted) adjacency matrix, and \mathbf{X} is the matrix of node signals/attributes. The adjacency matrix $\mathbf{A} \in \mathbb{R}^{N \times N}$ encodes the connection between two nodes (or joints), where if v_i, v_j are connected, then $A_{ij} = 1$, otherwise $A_{ij} = 0$. If each joint is endowed with a vectorized signal of 3D coordinates, *i.e.*, $\mathbf{x} : \mathcal{V} \rightarrow \mathbb{R}^3$, we may stack the signals of all nodes to form the signal matrix $\mathbf{X} \in \mathbb{R}^{N \times 3}$, where each row is associated with one node.

In order to extract skeletal graph features, we aim to perform convolutional filtering on these irregular attribute graphs like on regular grid-shape images. As studied in [33], [34], it is difficult to express a meaningful translation operator in the vertex domain. According to the spectral graph theory [31], the convolutional filtering on graphs depends on the graph Laplacian operator $\mathcal{L} = \mathbf{D} - \mathbf{A}$, where $\mathbf{D} \in \mathbb{R}^{N \times N}$ is the diagonal degree matrix with $D_{ii} = \sum_j A_{ij}$. For the graph Laplacian matrix, the normalized version is often used, *i.e.*,

$$\mathcal{L}^{norm} = \mathbf{D}^{-\frac{1}{2}} \mathcal{L} \mathbf{D}^{-\frac{1}{2}} = \mathbf{I} - \mathbf{D}^{-\frac{1}{2}} \mathbf{W} \mathbf{D}^{-\frac{1}{2}}, \quad (1)$$

where \mathbf{I} is the identity matrix. Unless otherwise specified, we use the normalized Laplacian matrix below.

As a real symmetric positive definite (SPD) matrix, the graph Laplacian matrix \mathcal{L} may be decomposed into

$$\mathcal{L} = \mathbf{U} \mathbf{\Lambda} \mathbf{U}^T, \quad (2)$$

where $\mathbf{\Lambda} = \text{diag}([\lambda_1, \lambda_2, \dots, \lambda_N])$ is a diagonal matrix of nonnegative real eigenvalues $\{\lambda_i\}$ (a.k.a. spectrum), and the

orthogonal matrix $\mathbf{U} = [\mathbf{u}_1, \dots, \mathbf{u}_N]$ means the corresponding eigenvectors. Analogous to the classic Fourier transform, the graph Fourier transform of a signal \mathbf{x} in spatial domain can be defined as $\hat{\mathbf{x}} = \mathbf{U}^\top \mathbf{x}$, where $\hat{\mathbf{x}}$ is the produced frequency signal. The corresponding inverse Fourier transform is $\mathbf{x} = \mathbf{U}\hat{\mathbf{x}}$.

Give any one filtering function $g(\cdot)$ of the graph \mathcal{L} , we can define the frequency responses on the input signal \mathbf{x} as $\hat{z}(\lambda_l) = \hat{x}(\lambda_l)\hat{g}(\lambda_l)$, or the inverse graph Fourier transform,

$$z(i) = \sum_{l=1}^N \hat{x}(\lambda_l)\hat{g}(\lambda_l)\hat{u}_l(i), \quad (3)$$

where $\hat{z}(\lambda_l), \hat{x}(\lambda_l), \hat{g}(\lambda_l)$ are the Fourier coefficients corresponding to the spectrum λ_l . Hence, the matrix description of graph filtering is

$$\mathbf{z} = \hat{g}(\mathcal{L})\mathbf{x} = \mathbf{U} \text{diag}([\hat{g}(\lambda_1), \dots, \hat{g}(\lambda_N)])\mathbf{U}^\top \mathbf{x}. \quad (4)$$

We need to learn the filtering function $g(\cdot)$, but the computational cost of Eqn. (2) and Eqn. (4) is expensive because the eigenvalue decomposition need to be done.

To address this problem, we may parameterize the filtering function $g(\cdot)$ with a polynomial approximation. As used in the literature [32], we employ the Chebyshev expansion of K order [35] by defining the recurrent relation $T_k(x) = 2xT_{k-1}(x) - T_{k-2}(x)$ with $T_0 = 1$ and $T_1 = x$. Any one function in the space $x \in [-1, 1]$ can be expressed with the expansion: $f(x) = \sum_{k=0}^{\infty} a_k T_k(x)$. Suppose we take the K -order approximation for $g(\cdot)$, *i.e.*,

$$\hat{g}(\lambda_l) = \sum_{k=0}^{K-1} \theta_k T_k(\tilde{\lambda}_l), \quad (5)$$

where $\theta = [\theta_1, \dots, \theta_K]^\top \in \mathbb{R}^K$ is the parameter vector of polynomial coefficients, and $\tilde{\lambda}_l = \frac{2}{\lambda_{max}}\lambda_l - 1$ with $\lambda_{max} = \max\{\lambda_1, \dots, \lambda_N\} \leq 2$. By substituting Eqn. (5) into Eqn. (4), we can derive the following equation,

$$\mathbf{z} = \sum_{k=0}^{K-1} \theta_k T_k\left(\frac{2}{\lambda_{max}}\mathcal{L} - \mathbf{I}\right)\mathbf{x}, \quad (6)$$

where we use this basic equation,

$$\begin{aligned} \mathcal{L}^k &= \mathbf{U} \text{diag}([\lambda_1^k, \dots, \lambda_N^k])\mathbf{U}^\top \\ &= (\mathbf{U} \text{diag}([\lambda_1, \dots, \lambda_N])\mathbf{U}^\top)^k. \end{aligned} \quad (7)$$

Let $\tilde{\mathcal{L}} = \frac{2}{\lambda_{max}}\mathcal{L} - \mathbf{I}$, if we denote the K filter bases about the Laplacian matrix as $\{T_0(\tilde{\mathcal{L}}), T_1(\tilde{\mathcal{L}}), \dots, T_{K-1}(\tilde{\mathcal{L}})\}$, the final local spectral graph filtering can be written as

$$\mathbf{Z} = [T_0(\tilde{\mathcal{L}})\mathbf{X}, T_1(\tilde{\mathcal{L}})\mathbf{X}, \dots, T_{K-1}(\tilde{\mathcal{L}})\mathbf{X}]\Theta, \quad (8)$$

where $\Theta \in \mathbb{R}^{3K \times d_z}$ is the parameters to be learnt with the d_z output channels, and $\mathbf{Z} \in \mathbb{R}^{N \times d_z}$ is the responses of spectral graph filtering. As $\tilde{\mathcal{L}}^k$ encodes a k -hop local neighborhood of each node, so the K -order polynomial in Eqn. (6) is a exactly K -localized filter function of graphs. Correspondingly, Θ is the K -localized filtering parameter need to be solved.

C. Detecting Salient Action Units

As observed that an action often occurs at a special body part, the detection of salient action units is necessary to reduce the disturbances of irrelevant joints as well as verify human cognition on actions. To detect action units, we propose a new layer named as action-attending layer to adaptively weight skeletal joints for different human actions, inspired by the attention mechanism used for various tasks, *e.g.*, machine translation [36], speech recognition [37], and image captioning [38], etc. Our purpose is to decide what action units identify (or play a key role in) a special human action.

We stack this layer after the spectral graph filtering layer to take advantage of high-level features. As the K -order filtering has a K -hop receptive field as introduced in Section III-B, the filtering responses of each node actually assemble certain information of K -path neighbors around the node. Suppose the feature responses after spectral filtering as $\mathbf{Z} \in \mathbb{R}^{N \times d_z}$ (in Eqn. 8), we attempt to learn a projecting matrix to weight nodes. But considering different action cases, we expect that the projecting matrix can dynamically change with the different input \mathbf{Z} . That is, the dynamic matrix $\mathcal{W} \in \mathbb{R}^{N' \times N}$ is actually a parameterized function of the variable \mathbf{Z} , formally,

$$\tilde{\mathbf{Z}} = \mathcal{W}(\mathbf{Z})\mathbf{Z}, \quad (9)$$

$$\begin{aligned} \text{s.t. } \mathcal{W}_{ij} &\geq 0, \quad \sum_{k=1}^N \mathcal{W}_{ik} = 1, \\ i &= 1, \dots, N', \quad j = 1, \dots, N'. \end{aligned} \quad (10)$$

The larger the matrix item \mathcal{W}_{ij} is, the more important the j -the node is for action recognition. To model the dynamic property, we define the dynamic function as

$$\mathcal{W}(\mathbf{Z}) = (\tanh(\mathbf{Z}\mathbf{Q} + \mathbf{b}^\top)\mathbf{V})^\top, \quad (11)$$

$$\mathcal{W}_{ij} = \frac{\exp(\mathcal{W}_{ij})}{\sum_{k=1}^N \exp(\mathcal{W}_{ik})}, \quad (12)$$

where $\mathbf{Q} \in \mathbb{R}^{d_z \times d'}$, $\mathbf{V} \in \mathbb{R}^{d' \times N'}$, $\mathbf{b} \in \mathbb{R}^{d'}$ are the parameters to be solved. Hence, the dynamic function \mathcal{W} takes advantage of the input feature \mathbf{Z} .

In addition to the detection of salient action units, the dynamic weighting function has two extra advantages:

(1) Order-independency of nodes.

Suppose all nodes are ordered in $\{v_1, \dots, v_N\}$ and the signal matrix $\mathbf{X} = [\mathbf{x}_1, \dots, \mathbf{x}_N]^\top$ is built, then we can extract the filtering feature $\mathbf{Z} = [\mathbf{z}_1, \dots, \mathbf{z}_N]^\top$ in the same order of nodes, *i.e.*, \mathbf{z}_i is still associated with v_i . However, if all nodes are disordered, *e.g.*, $\{v_N, v_{N-1}, \dots, v_1\}$, correspondingly, we have $\mathbf{X}' = [\mathbf{x}_N, \mathbf{x}_{N-1}, \dots, \mathbf{x}_1]^\top$ and $\mathbf{Z}' = [\mathbf{z}_N, \mathbf{z}_{N-1}, \dots, \mathbf{z}_1]^\top$. If we employ a constant projection \mathbf{W} rather than the dynamic function \mathcal{W} , it will result into $\mathbf{W}\mathbf{Z} \neq \mathbf{W}\mathbf{Z}'$. That means, different traversing ways on a graph will produce different responses, which is unfeasible to feature comparisons, *e.g.*, the feature $\tilde{\mathbf{Z}}$ will be flatten as the input to fed into recurrent neural network (See Section III-D). However, the dynamic weighting function \mathcal{W} is order-independent for graphical nodes, according to the following theory.

Theorem 1. Given the dynamic function \mathcal{W} in Eqn. (11), the output $\tilde{\mathbf{Z}}$ in Eqn. (9) is irrelevant to the order of traversing order of graphical nodes.

Proof. Suppose the input signal matrix \mathbf{X} , correspondingly, we can obtain the convolutional filtering feature \mathbf{Z} and $\tilde{\mathbf{Z}}$ according to Eqn. (8) and Eqn. (9) respectively. Given any one permutation matrix $\mathcal{P} \in \{0, 1\}^{N \times N}$, $\mathcal{P}\mathbf{X}$ equals to reorder all signals in \mathbf{X} . Let denote the feature \mathbf{Z}' and $\tilde{\mathbf{Z}}'$ when taking $\mathcal{P}\mathbf{X}$ as the input. According to Eqn. (8), we can have

$$\begin{aligned} \mathbf{Z}' &= [T_0(\mathcal{P}\tilde{\mathcal{L}}\mathcal{P}^\top)\mathcal{P}\mathbf{X}, \dots, T_{K-1}(\mathcal{P}\tilde{\mathcal{L}}\mathcal{P}^\top)\mathcal{P}\mathbf{X}]\Theta \\ &= [\mathcal{P}T_0(\tilde{\mathcal{L}})\mathcal{P}^\top\mathcal{P}\mathbf{X}, \dots, \mathcal{P}T_{K-1}(\tilde{\mathcal{L}})\mathcal{P}^\top\mathcal{P}\mathbf{X}]\Theta \\ &= \mathcal{P}[T_0(\tilde{\mathcal{L}})\mathbf{X}, \dots, T_{K-1}(\tilde{\mathcal{L}})\mathbf{X}]\Theta \\ &= \mathcal{P}\mathbf{Z}. \end{aligned} \quad (13)$$

Note that the Laplacian matrix is also reordered by the permutation matrix \mathcal{P} . Now we only need to prove that $\tilde{\mathbf{Z}} = \tilde{\mathbf{Z}}'$. According to Eqn. (9) and Eqn. (11), we have

$$\begin{aligned} \tilde{\mathbf{Z}}' &= \mathcal{W}(\mathbf{Z}')\mathbf{Z}' = \mathcal{W}(\mathcal{P}\mathbf{Z})\mathcal{P}\mathbf{Z} \\ &= (\tanh(\mathcal{P}\mathbf{Z}\mathbf{Q} + \mathbf{b}^\top)\mathbf{V})^\top\mathcal{P}\mathbf{Z} \\ &= (\mathcal{P}\tanh(\mathbf{Z}\mathbf{Q} + \mathbf{b}^\top)\mathbf{V})^\top\mathcal{P}\mathbf{Z} \\ &= (\tanh(\mathbf{Z}\mathbf{Q} + \mathbf{b}^\top)\mathbf{V})^\top\mathcal{P}^\top\mathcal{P}\mathbf{Z} \\ &= (\tanh(\mathbf{Z}\mathbf{Q} + \mathbf{b}^\top)\mathbf{V})^\top\mathbf{Z} \\ &= \tilde{\mathbf{Z}}. \end{aligned} \quad (14)$$

□

(2) Dynamic pooling of nodes.

The dynamic function \mathcal{W} may be regarded as a pooling operation on nodes. Given a row \mathcal{W}_i of \mathcal{W} , the output $\mathcal{W}_i\mathbf{Z}$ is a new signal by weighting and combining nodes. Correspondingly, we can update the adjacency matrix of nodes and the Laplacian matrix,

$$\mathbf{A}' = \mathcal{W}\mathbf{A}\mathcal{W}^\top, \quad (15)$$

$$\tilde{\mathcal{L}}' = \mathbf{I} - \mathbf{D}^{-1/2}\mathbf{A}'\mathbf{D}^{-1/2}, \quad (16)$$

where $\mathbf{A}', \tilde{\mathcal{L}}' \in \mathbb{R}^{N' \times N'}$. Thus, the new Laplacian matrix $\tilde{\mathcal{L}}'$ and the feature $\tilde{\mathbf{Z}}$ can be fed into the next layer and make the neural network go deeper.

D. Modeling Temporal Motions

After extracting the spatial graphical features, we need to model temporal motion variations of a skeletal sequence. Many non-linear dynamic models can be used to solve this problem. Here we employ a special class of recurrent neural networks (RNN), long short-term memory (LSTM) [39], which can mitigate gradient vanishment when back-propagating gradients. LSTM has demonstrated the powerful ability to model long-range dependencies [40], [41], [42]. Suppose the graphical feature of the t -th skeletal frame is $\tilde{\mathbf{z}}_t = \text{vectorize}(\tilde{\mathbf{Z}}_t) \in \mathbb{R}^{N'd_z}$, the cell output $\mathbf{h}_t \in \mathbb{R}^{d_h}$ and

states $\mathbf{c}_t \in \mathbb{R}^{d_h}$ are intermediate vectors, formally, the motion variations can be modeled as

$$\mathbf{i} = \sigma(\mathbf{W}_{zi}\tilde{\mathbf{z}}_t + \mathbf{W}_{hi}\mathbf{h}_{t-1} + \mathbf{w}_{ci} \odot \mathbf{c}_{t-1} + \mathbf{b}_i), \quad (17)$$

$$\mathbf{f} = \sigma(\mathbf{W}_{zf}\tilde{\mathbf{z}}_t + \mathbf{W}_{hf}\mathbf{h}_{t-1} + \mathbf{w}_{cf} \odot \mathbf{c}_{t-1} + \mathbf{b}_f), \quad (18)$$

$$\mathbf{c}_t = \mathbf{f}_t \odot \mathbf{c}_{t-1} + \mathbf{i}_t \odot \tanh(\mathbf{W}_{zc}\tilde{\mathbf{z}}_t + \mathbf{W}_{hc}\mathbf{h}_{t-1} + \mathbf{b}_c), \quad (19)$$

$$\mathbf{o}_t = \sigma(\mathbf{W}_{zo}\tilde{\mathbf{z}}_t + \mathbf{W}_{ho}\mathbf{h}_{t-1} + \mathbf{w}_{co} \odot \mathbf{c}_t + \mathbf{b}_o), \quad (20)$$

$$\mathbf{h}_t = \mathbf{o}_t \odot \tanh(\mathbf{c}_t), \quad (21)$$

where $\sigma(\cdot)$ is the elementwise sigmoid function, *i.e.*, $\sigma(\mathbf{x}) = 1/(1 + e^{-\mathbf{x}})$, \odot denotes the Hadamard product and $\mathbf{i}, \mathbf{f}, \mathbf{o}, \mathbf{c} \in \mathbb{R}^{d_h}$ are respectively the *input gate*, *forget gate*, *output gate*, *cell* and *cell input* activation vectors. The weight matrices $\{\mathbf{W}_{z\cdot} \in \mathbb{R}^{d_h \times N'd_z}, \mathbf{W}_{h\cdot} \in \mathbb{R}^{d_h \times d_h}, \mathbf{w}_c \in \mathbb{R}^{d_h}\}$ and the bias vectors $\{\mathbf{b}_i, \mathbf{b}_f, \mathbf{b}_c, \mathbf{b}_o \in \mathbb{R}^{d_h}\}$ are the model parameters to be solved. Finally, we use the output gate \mathbf{o}_t as the response at the t -th time slice.

IV. IMPLEMENTATION DETAILS

In this section, we will introduce our implementation details including network architecture and data augmentation.

A. Network Architecture

For the undirected graph, we simply construct edge connections based on human bones. That means, if two joints are bridged with a bone, the edge weight is assigned to 1, otherwise 0. Our plain network contains two spectral graph filtering layers along with two action-attending layers followed by a temporal recurrent layer, as shown in Fig. (1). In the spectral filtering layers, the receptive fields are set to $K = 10$ as default, and the output signals have the length of 32 dimensions and 64 dimensions respectively for two layers. In the action-attending layer, we take a simple design rule: the dimensions remain invariant with regard to the input, *i.e.*, $N' = N, d' = d_z$. In the temporal recurrent layer, we employ the classic LSTM unit to model the temporal dynamics, where the dimension of hidden units is set to 256. In the full connected layer, the output has the same dimension to the input. The network ends with a cross-entropy loss used for classification. The learning rate of network is set to 0.02 with a momentum of 0.9. More analysis/discussion of parameters can be found in Section V-E. The concrete implementation takes TensorFlow as the infrastructure.

B. Data Processing

As skeletal data is usually captured from multi-view points and human actions are independent on the user coordinate system, we modify the origin of the coordinate system as the orthocenter of joints for each frame of skeleton, *i.e.*, $\mathcal{O} = \frac{1}{N} \sum_{i=1}^N \mathbf{x}_i$, where $\mathbf{x}_i \in \mathbb{R}^3$ is a 3D coordinate of the i -th joint, N is the number of joints.

To enhance the robustness of model training, we perform data augmentation as widely used in previous deep learning literature [43], [20]. Concretely, for each action sequence, we split the sequence into several equal sized subsequences, here

TABLE I
SUMMARIZATION OF FOUR ACTION RECOGNITION DATASETS.

Dataset	# Joints	# Actions	# Subjects	# Sequences
HDM05 [17]	31	130	5	2337
Florence 3D [18]	15	9	10	215
LSC [19]	15/20	88	79	3898
NTU RGB+D [20]	25	60	40	56880

12 segments, and then pick one frame from each segment randomly to generate a large amount of training sequences. In addition, we randomly scale the skeletons by multiplying a factor in [0.98, 1.02] for the sake of the adaptive capability of scaling.

V. EXPERIMENTS

To evaluate our proposed A²GNN, we conduct extensive experiments on four benchmark skeleton-based action datasets, including HDM05 [17], Florence 3D [18], Large Scale Combined dataset (LSC) [19] and NTU RGB+D [20]. A brief summarization about them is given in Table I. Below we will compare our A²GNN with the recent state-of-the-art methods, then analyze confusion matrices and action unit detection, and finally discuss some network parameters.

A. Datasets

1) *HDM05 [17]*: This dataset was captured by using an optical marker-based Vicon system, and gathered 2337 action sequences for 130 motion classes, which are performed by 5 non-professional actors named “bd”, “bk”, “dg”, “mm” and “tr”. Each skeleton data is represented with 31 joints. Until now, this dataset should involve the most skeleton-based action categories to the best of our knowledge. Due to the intra-class variations and large number of motion classes, this dataset is challenging in action recognition.

2) *Florence 3D [18]*: This dataset was collected via a stationary Microsoft Kinect camera. It consists of 215 action sequences from 10 subjects for 9 actions: wave, drink from a bottle, answer phone, clap, tight lace, sit down, stand up, read watch, bow. Only 15 joints are recorded for each skeletal data. As a few skeletal joints, some types of actions are difficult to distinguish, such as drink from a bottle, answer phone and read watch.

3) *LSC [19]*: This dataset combines nine publicly available datasets, including MSR Action3D Ext [44], [45], [46], UTKinect-Action3D [12], MSR DailyActivity 3D [27], MSR Action Pairs 3D [47], CAD120 [48], CAD60 [49], G3D [50], [51], RGBD-HuDa [52], UTD-MHAD [53], and form a complex action dataset with 94 actions. As some samples have not skeleton information, we remove them and construct a skeleton dataset of 88 actions by following previous standard protocols. As each individual dataset has its own characteristics in the action execution manners, backgrounds, acting positions, view angles, resolutions, and sensor types, the combination of a large number of action classes makes the dataset more challenging in suffering large intra-class variation compared to each individual dataset.

TABLE II
COMPARISONS ON HDM05 DATASET.

Method	Protocol [54]	Protocol [55]
	Accuracy	Accuracy
RSR-ML [56]	40.0%	-
Cov-RP [57]	58.9%	-
Ker-RP [54]	66.2%	-
SPDNet [55]	-	61.45%±1.12
Lie Group [10]	-	70.26%±2.89
LieNet [58]	-	75.78%±2.26
P-LSTM [20]	70.4%	73.42%±2.05
A²GNN	76.5%	84.47%±1.52

*Note that all 130 classes are used here.

TABLE III
COMPARISONS ON FLORENCE DATASET.

Method	Accuracy
Multi-part Bag-of-Poses [18]	82.00%
Riemannian Manifold [6]	87.04%
Lie Group [10]	90.88%
Graph-Based [13]	91.63%
MIMTL [59]	95.29%
P-LSTM [20]	95.35%
A²GNN	98.60%

4) *NTU RGB+D [20]*: This dataset is collected by Microsoft Kinect v2 cameras from different views. It consists of 56880 sequences and over 4 million frames for 60 distinct actions, including various of daily actions and pair actions. These actions were performed by 40 subjects aged between 10 and 35. The skeleton data is represented by 25 joints. As far as we know, this dataset is currently the largest skeleton-based action recognition dataset. The large intra-class and view point variations make this dataset great challenging. Meanwhile, a large amount of samples will bring a new challenging to the current skeleton-based action recognition methods.

B. Comparisons with State-of-the-art Methods

1) The Results:

a) *HDM05*: To compare with those previous literatures, we conduct two types of experiments by following two widely-used protocols. First, we follow the protocol used in [54] to perform action recognition on all of the 130 classes. Actions of two subjects named “bd” and “mm” are used to train the model, and the remaining three ones for testing. The comparison results are shown in the left column of Table II. Second, to fairly compare the current deep learning methods, we follow the settings of the literature [55] to conduct 10 random evaluations, each of which randomly selects half of the sequences for training and the rest for testing. The results are reported in the right column of Table II.

b) *Florence 3D*: We follow the experimental settings of the literature [13] to perform leave-one-subject-out cross-validation. In each round, skeletal data of 9 subjects is taken for training and the remain one for testing. The experimental results are reported in Table III.

TABLE IV
COMPARISONS ON LARGE SCALE COMBINED DATASET.

Method	Cross Sample		Cross Subject	
	Precision	Recall	Precision	Recall
HON4D [47]	84.6%	84.1%	63.1%	59.3%
Dynamic Skeletons [60]	85.9%	85.6%	74.5%	73.7%
P-LSTM [20]	84.2%	84.9%	76.3%	74.6%
A²GNN	87.6%	88.1%	84.0%	82.0%

TABLE V
COMPARISONS ON NTU RGB+D DATASET.

Method	Cross Subject	Cross View
	Accuracy	Accuracy
HON4D [47]	30.56%	7.26%
Lie Group [10]	50.08%	52.76%
Skeletal Quads [61]	38.62%	41.36%
Dynamic Skeletons [60]	60.23%	65.22%
HBRNN [62]	59.07%	63.97%
LieNet [58]	61.37%	66.95%
Deep RNN [20]	56.29%	64.09%
Deep LSTM [20]	60.69%	67.29%
P-LSTM [20]	62.93%	70.27%
ST-LSTM [43]	69.2%	77.7%
A²GNN	72.74%	82.80%

c) *LSC*: We follow the protocol designed in the recent work [19] to conduct two types of experiments, random cross subject and random cross sample. The experimental results are reported in Table IV.

d) *NTU RGB+D*: Different from the above three datasets, we preprocess the joint coordinates in a way similar to [20]. Concretely, after translating the original coordinates of body joints as mentioned above, we rotate the x axis parallel to the 3D vector from “right shoulder” to “left shoulder” and y axis towards the 3D vector from “spine base” to “spine”. The z axis is fixed as the new $x \times y$. This dataset has two types of standard evaluation protocols [20]. One is cross-subject evaluation, for which half of the subjects are used for training and the remaining ones for testing. The second is cross-view evaluation, for which two viewpoints are used for training and one is left out for testing. The experimental results are shown in Table V.

2) The Analysis:

As shown in Table II~V, we compare the current state-of-the-art methods on different datasets, including the shallow learning methods and the deep learning methods. From these results, we have the following observations:

- *Matrix-based descriptors* (e.g., *covariance or its variants*) are conventionally used to model spatial or temporal relationships. Moreover, these descriptors are often regarded to be embedded on specific geometric manifolds [10], [58], [55]. The advanced variant [54] improved the performance by constructing some robust SPD matrices. More recently, SPDNet [58] and LieNet [55] attempted to learn deep features from those raw matrix descriptors under the assumption of manifold. Although the deep manifold learning strategy on matrix descriptors raises a promising

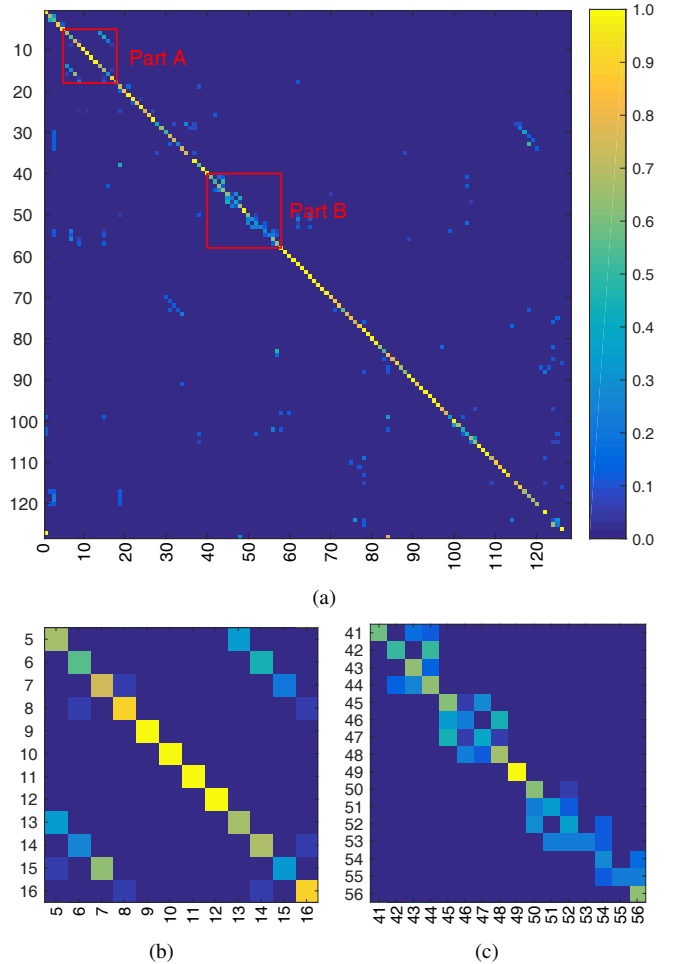


Fig. 2. Confusion matrix of our A²GNN on HDM05 dataset according to the testing protocol in [55]. (b) and (c) are the close up of Part A and Part B in (a). X-axis and Y-axis are associated with the indices of action classes. ** Indices of part classes: 5-depositFloorR; 6-depositHighR; 7-depositLowR; 8-depositMiddleR; 13-grabFloorR; 14-grabHighR; 15-grabLowR; 16-grabMiddleR; 41-kickLFront1Reps; 42-kickLFront2Reps; 43-kickLSide1Reps; 44-kickLSide2Reps; 45-kickRFront1Reps; 46-kickRFront2Reps; 47-kickRSide1Reps; 48-kickRSide2Reps; 50-punchLFront1Reps; 51-punchLFront2Reps; 52-punchLSide1Reps; 53-punchLSide2Reps; 54-punchRFront1Reps; 55-punchRFront2Reps; 56-punchRSide1Reps.

direction to some extent, the matrix-based representations principally limit their capability of modeling dynamic variations, because the only second-order statistic relationship of skeletal joints is preserved in the descriptors, whereas first-order statistics is also informative [63].

- *Deep features are more effective than those shallow features for skeleton-based action representation.* The advanced nonlinear dynamic networks, specifically DeepRNN, Deep-LSTM, P-LSTM [20] and ST-LSTM [43], largely improve the action recognition performance, due to the good encoding capability of gated network units. Most of them use recurrent networks to model temporal dynamics. Besides, ST-LSTM also attempted to model spatial skeletal joints by taking a tree-structure traversal way on spatial joints. Similar to them, we also use recurrent neural network to model temporal dynamics. But different from them, we directly extract high-level

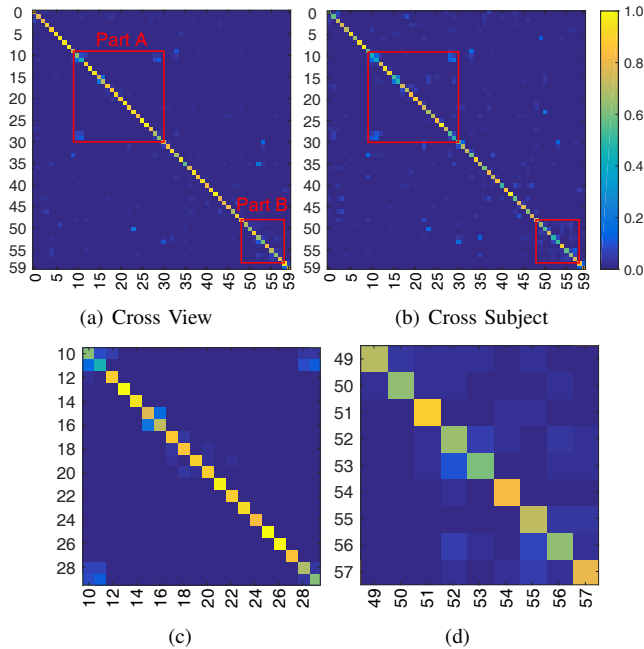


Fig. 3. Confusion matrix of our A^2GNN on NTU RGB+D dataset. (c) and (d) are the close up of Part A and Part B in (a). X-axis and Y-axis are associated with the indices of action classes.

** Indices of part classes: 10-reading; 11-writing; 15-wear a shoe; 16-take off a shoe; 17-wear on glasses; 18-take off glasses; 19-put on a hat/cap; 20-take off a hat/cap; 28-playing with phone/tablet; 29-typing on a keyboard; 49-punching/slapping other person; 50-kicking other person; 51-pushing other person; 52-pat on back of other person; 53-point finger at the other person; 54-hugging other person; 55-giving something to other person; 56-touch other person's pocket; 57-handshaking.

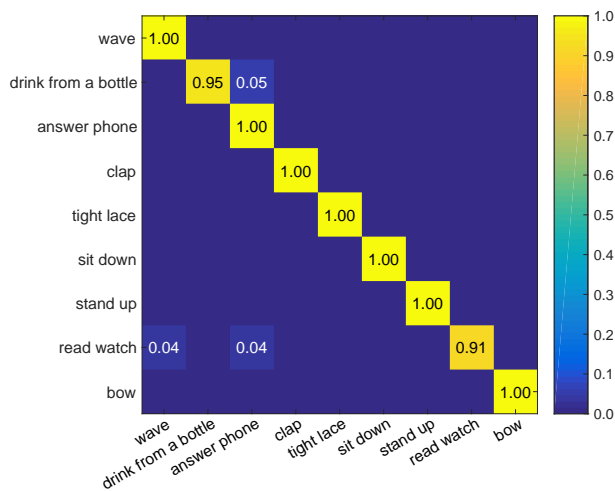


Fig. 4. Confusion matrix of our A^2GNN on Florence 3D Actions dataset.

semantic features from spatial skeletal graphs like the standard convolutional neural network.

- *The proposed deep graph method is superior to the recent graph-based method [13].* As shown in Table III, our A^2GNN has a large improvement (about 7%) in contrast to the work [13]. In principle, our A^2GNN is very different from this work [13], although graph is used for both. First, the graph is used to describe skeleton at each temporal slice for our method, not those segmented

motion parts (called motionlets). Second, the purpose of the use of graph is to extract skeletal features for ours, not model the relationship of motion parts. Third, our method is a fully end-to-end deep learning architecture, not abide by the conventional two-step way: i) construct graphs of motion parts, and ii) compute the similarity between graphs via subgraph-pattern graph kernel. Besides, the action-attending mechanism is introduced into our network architecture.

- *Our proposed A^2GNN greatly improves the current state-of-the-art on most datasets.* On the Florence dataset, our method achieves a nearly perfect performance 98.60%. On the current largest dataset NTU RGB+D dataset, the state-of-the-art performance is pushed to the higher 72.74% and 82.80%, from 69.2% and 77.7% for ST-LSTM. In summary, we can benefit from the deep graph network architecture. The reason can be two folds: i) the deep skeletal graph features, rather than simple spatial or temporal features learnt by LSTM; and ii) the preservation of more original feature information (as signals of each node), rather than only the second-order statistics of skeletal joints.
- *Different performances occur on different datasets.* Among the four datasets, Florence 3D is the simplest one with 215 sequences and 9 action classes, so most methods can obtain a good accuracy. The most difficult dataset should be the largest dataset NTU RGB+D dataset, which consists of 56880 sequences and covers various of daily actions and pair actions. The cross subject accuracy on it just surpasses 70% due to various entangled actions as analyzed in Section V-C. Specifically, the phenomenon of entangled actions deteriorates in HDM05, which contains many confusable classes, such as walk step number, walk start with left or right, etc.
- *Cross subject is more difficult than cross view or cross sample.* The phenomenon is observed from Table IV, Table V, and Table II (the left/right column w.r.t cross subject/cross sample). It is easy to understand, each subject has itself action characteristics. In the cross subject task, more unforeseeable information exists the testing set, compared to the other tasks.

C. Analysis of Confusion Matrices

To further reveal what classes are easy to confuse with others, we show confusion matrices on HDM05, NTU RGB+D and Florence datasets, respectively in Fig. 2, Fig. 3, and Fig. 4. For LSC dataset, we don't depict its confusion matrix due to the different evaluation criterion (precision and recall).

For the HDM05 dataset, as shown in Fig. 2(a), we give the confusion matrix of 130 classes in the case of recognition rate 84.47% (i.e., the testing protocol follows the literature [55]). The diagonal characterizes the correct classification for each action, and those non-diagonals depict the confusion results cross different classes. In order to view them more clearly, we provide two close up areas (Part A and Part B) in Fig. 2(b) and Fig. 2(c). As observed from the two subfigures, our proposed method suffers some failures in

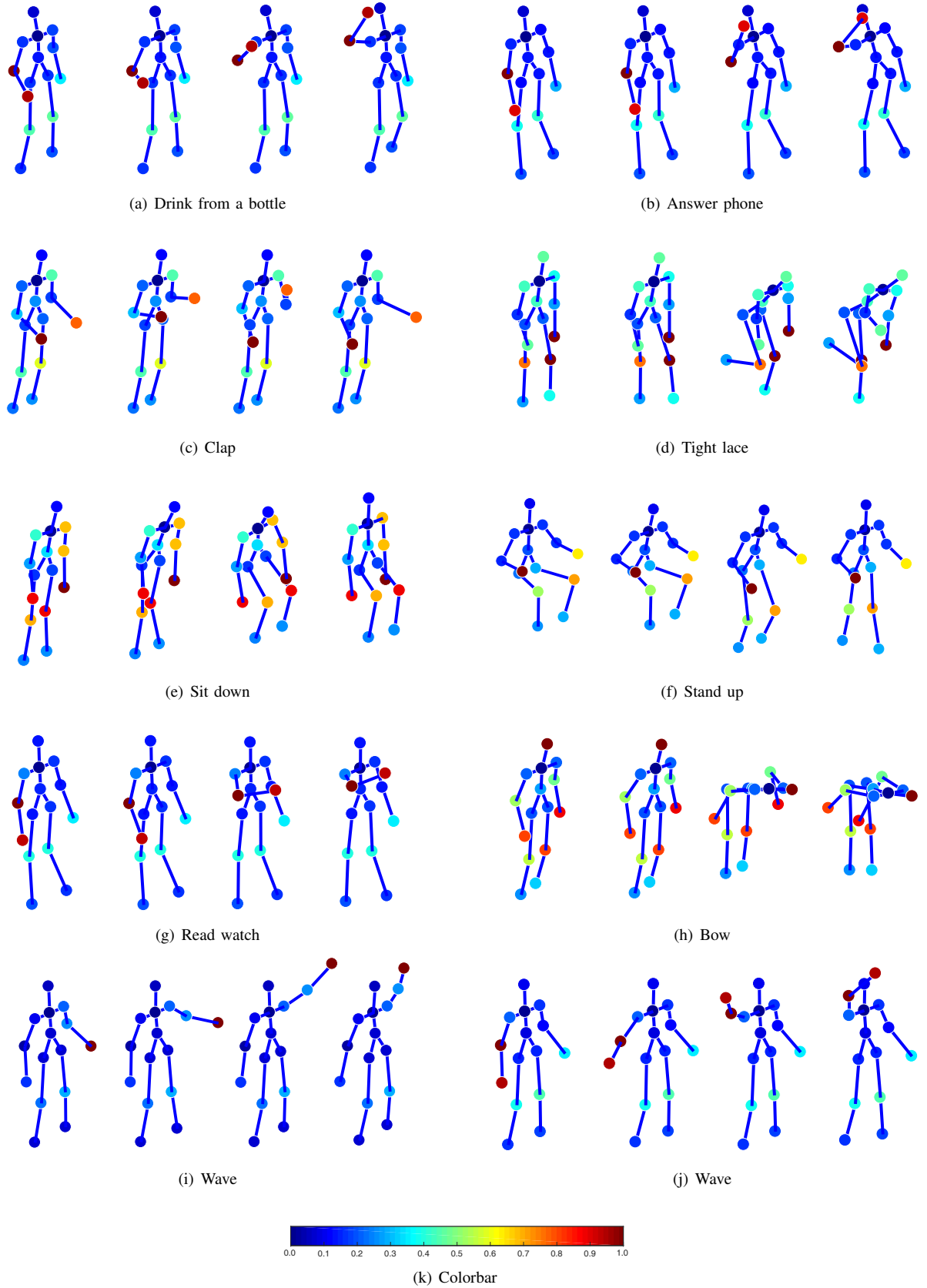


Fig. 5. Visual examples of detected salient action units of all 9 different action classes on Florence 3D. Higher weights in the colorbar means more important to characterize the action. Note that the motion sequences in (i) and (j) are annotated as one class (“wave”) in this dataset, although one is left arm while the other is right arm.

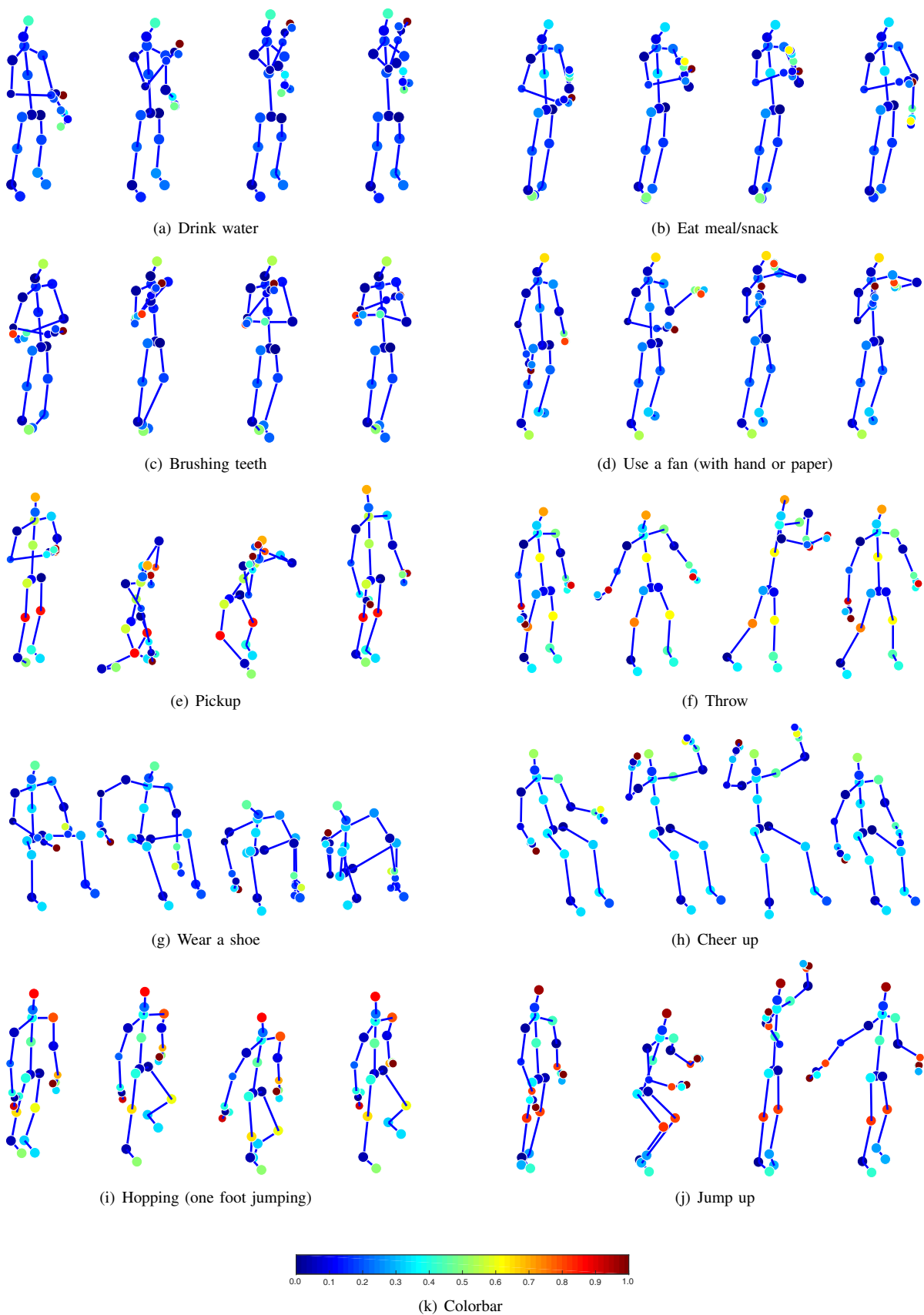


Fig. 6. Visual examples of detected salient action units on NTU RGB+D. Higher weights in the colorbar means more important to characterize the action.

TABLE VI
RESULTS OF TUNING NETWORK MODULES.

Method	HDM05		Florence	LSC				NTU RGB+D	
	Protocol [54]	Protocol [55]		Cross Sample		Cross Subject		Cross Subject	Cross View
	Acc.	Acc.		Acc.	Prec.	Rec.	Prec.	Rec.	Acc.
A ² GNN \ominus filt. \ominus AU det.	71.67%	77.74% \pm 1.61	93.49%	79.03%	76.96%	73.21%	70.77%	63.18%	68.82%
A ² GNN \ominus AU det.	75.28%	83.01% \pm 1.43	96.74%	86.56%	85.71%	79.97%	79.92%	68.26%	80.08%
A ² GNN	76.46%	84.47% \pm 1.52	98.60%	87.61%	88.10%	83.98%	82.03%	72.74%	82.80%

distinguishing those very similar actions. For example, the depositing and grabbing are seriously confused, as “deposit-FloorR” and “grabFloorR” are almost visually consistent in knees bent and arms stretch. For a basic kicking behavior, the actions of left foot forwards (“kickLFront1Reps”) or left foot sideways (“kickLSide1Reps”) are subtle with small inter-class distance. Likewise, there are some other confusable actions, such as “kickRFront1Reps” vs “kickRSide1Reps”, “punchLFront1Reps” vs “punchLFront2Reps” and so on. Note that the postfixes “1Reps” and “2Reps” means the repetitive times of a action.

Different from HDM05 with finer-partitioned action classes, NTU RGB+D contains many pairs of reversed actions, such as “wear A” vs “take off A”, “put on A” vs “take off A”, etc. As shown in Fig. 3(c), these reversed pairs are easy to be confused when only considering human skeleton data. Besides, similar to the observation from HDM05, the confusion cases occur in those similar motions, such as “pat on back of other person” vs “point finger at the other person”. Likewise, this phenomenon happens in Florence 3D, as shown in Fig. 4. Although the accuracies of 100% are achieved on seven actions (*i.e.*, wave, answer phone, clap, tight lace, sit down, stand up, bow), there are a few uncorrect classifications including “drink from a bottle” to “answer phone”, “read watch” to “wave”/“answer phone”. A future possible solution to these above confusable phenomena is that the contextual information in the RGB color space may be considered to compensate for the plain coordinate information of skeletal joints to some extent.

D. Visualization of Action Units

To verify whether the action-attending layer detects those salient action units, we provide some visualization examples in Fig. 5 and Fig. 6, where skeletal joints are colored according to the learnt weights. Note that here we exhibit the first action-attending layer because the learnt weights are directly associated with skeletal joints. From the visualization of Florence 3D in Fig. 5, we can find that our A²GNN is able to learn those salient joints for all 9 actions. For examples, for the “wave” action, those joints on the moving right arm are important; for the “tight lace” action, the joints of hands and knees (bent) is endowed with higher weights. Specifically, for the same “wave” action, our method can still correctly detect the moving units of the left arm or the right arm, as shown in Fig. 5(i) and Fig. 5(j).

For the more complex action dataset, NTU RGB+D, we can still observe that those detected salient action units are almost matched to our intuitive understanding. For the “pick

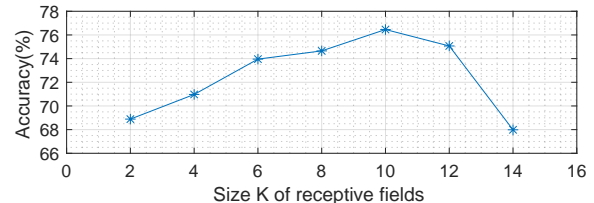


Fig. 7. The performance trend of different receptive fields K on HDM05.

up” action, the units of head, hands and knees are crucial to identify this action. For the “hopping” action (right foot jumping) in Fig. 6(i), besides those salient joints on head, hands and knees, the joint on left shoulder is still more salient than that of right shoulder due to the more drastic motion variations on left shoulder compared to right shoulder. In contrast to Florence 3D, as NTU RGB+D is configured with more sensors to record human motions, thus more detailed motions can be captured, *e.g.*, the toe joint on the “wear a shoe” action, the “hopping” action.

Consequently, the above observations indicate that the action-attending layer can adaptively weight skeletal joints for different actions, and the detected salient action units almost conform to our cognitive understanding on human actions. Moreover, the detection of action units can bring some gains of the accuracy as analyzed in Section V-E.

E. Tuning of Our Network

In our proposed A²GNN, there are three main modules: spectral graph filtering, action unit detection, and temporal motion modeling. The last module is an indispensable unit to our network. To dissect our network architecture, we conduct experiments on the four datasets by removing the former two modules (*i.e.*, A²GNN \ominus filt. \ominus AU det.) or only removing AU detection (*i.e.*, A²GNN \ominus AU det.). The comparison results are reported in Table VI. As observed from this table, the detection of action units can be helpful for action recognition more or less, as skeletal joints are successfully activated with different weights as analyzed in Section V-D. Further, the spectral graph filtering module plays a tremendous role in promoting recognition accuracy, as exhibited at the former two rows in the table. The main reason should be that high-level semantic features are extracted from graphs like the conventional convolution network on grid-shape images.

Among the parameters of our network, the most key is the size K of receptive fields. With the increase of K , the local filtering region, *i.e.*, covering the hopping neighbors,

will become larger. To check the influence of different K values, we conduct an experiment on HDM05 dataset by testing $K = 2, 4, 6, 8, 10, 12, 14$. The results are shown in Fig. 7, we can find that, the performance improves with the increase of K due to larger receptive fields, and then degrades after $K = 10$, for which a possible reason is the locality of salient action units.

VI. CONCLUSION

In this paper, an end-to-end action-attending graphic neural network (A^2GNN) is proposed to deal with the task of skeleton-based action recognition. In order to extract deep features from body skeletons, we model human skeletons into undirected attribute graphs, and then perform spectral graph filtering on skeletal graphs like the standard CNN. To detect those salient action units crucial to identify human motions, we further design an action-attending layer to adaptively weight skeletal joints for different actions. Those extracted deep graph features at consecutive frames are finally fed into a recurrent network of LSTM. Extensive experiments and analyses have indicated that the modules of spectral graph filtering and action unit detection play an important role in the improvement on action classification. Especially, the action-attending layer also produces some interesting salient action units, which may be understandable from the view of our cognition. Further, our proposed A^2GNN has achieved the state-of-the-art results on the four public skeleton-based action datasets, including the current largest and most challenging NTU RGB+D dataset. In the future, we will explore the fusion RGB color information into our network.

REFERENCES

- [1] J. K. Aggarwal and M. S. Ryoo, "Human activity analysis: A review," *ACM Comput. Survey*, vol. 43, no. 3, p. 16, 2011.
- [2] M. Ye, Q. Zhang, L. Wang, J. Zhu, R. Yang, and J. Gall, "A survey on human motion analysis from depth data," in *Time-of-Flight and Depth Imaging. Sensors, Algorithms, and Applications*, 2013, pp. 149–187.
- [3] G. Johansson, "Visual motion perception," *Scientific Am.*, vol. 232, no. 6, pp. 76–88, 1975.
- [4] B. B. Amor, J. Su, and A. Srivastava, "Action recognition using rate-invariant analysis of skeletal shape trajectories," *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 38, no. 1, pp. 1–13, 2016.
- [5] X. Cai, W. Zhou, L. Wu, J. Luo, and H. Li, "Effective active skeleton representation for low latency human action recognition," *IEEE Trans. Multimedia*, vol. 18, no. 2, pp. 141–154, 2016.
- [6] M. Devanne, H. Wannous, S. Berretti, P. Pala, M. Daoudi, and A. Del Bimbo, "3-d human action recognition by shape analysis of motion trajectories on riemannian manifold," *IEEE Trans. Cybern.*, vol. 45, no. 7, pp. 1340–1352, 2015.
- [7] C. Ellis, S. Z. Masood, M. F. Tappen, J. J. LaViola, and R. Sukthankar, "Exploring the trade-off between accuracy and observational latency in action recognition," *Int. J. Comput. Vis.*, vol. 101, no. 3, pp. 420–436, 2013.
- [8] R. Slama, H. Wannous, M. Daoudi, and A. Srivastava, "Accurate 3d action recognition using learning on the grassmann manifold," *Pattern Recognit.*, vol. 48, no. 2, pp. 556–567, 2015.
- [9] L. Tao and R. Vidal, "Moving poselets: A discriminative and interpretable skeletal motion representation for action recognition," in *Proc. ICCVW*, 2015, pp. 61–69.
- [10] R. Vemulapalli, F. Arrate, and R. Chellappa, "Human action recognition by representing 3d skeletons as points in a lie group," in *Proc. CVPR*, 2014, pp. 588–595.
- [11] J. Wang, Z. Liu, Y. Wu, and J. Yuan, "Learning actionlet ensemble for 3d human action recognition," *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 5, no. 36, pp. 914–927, 2014.
- [12] L. Xia, C.-C. Chen, and J. Aggarwal, "View invariant human action recognition using histograms of 3d joints," in *Proc. CVPRW*, 2012, pp. 20–27.
- [13] P. Wang, C. Yuan, W. Hu, B. Li, and Y. Zhang, "Graph based skeleton motion representation and similarity measurement for action recognition," in *Proc. ECCV*, 2016, pp. 370–385.
- [14] J. Shotton, A. Fitzgibbon, M. Cook, T. Sharp, M. Finocchio, R. Moore, A. Kipman, and A. Blake, "Real-time human pose recognition in parts from single depth images," in *Proc. CVPR*, 2011, pp. 1297–1304.
- [15] Y. Du, Y. Fu, and L. Wang, "Representation learning of temporal dynamics for skeleton-based action recognition," *IEEE Trans. Image Process.*, vol. 25, no. 7, pp. 3010–3022, 2016.
- [16] A. Krizhevsky, I. Sutskever, and G. E. Hinton, "Imagenet classification with deep convolutional neural networks," in *Proc. NIPS*, 2012, pp. 1097–1105.
- [17] M. Müller, T. Röder, M. Clausen, B. Eberhardt, B. Krüger, and A. Weber, "Documentation mocap database hdm05," 2007.
- [18] L. Seidenari, V. Varano, S. Berretti, A. Bimbo, and P. Pala, "Recognizing actions from depth cameras as weakly aligned multi-part bag-of-poses," in *Proc. CVPRW*, 2013, pp. 479–485.
- [19] J. Zhang, W. Li, P. Wang, P. Ogunbona, S. Liu, and C. Tang, "A large scale rgb-d dataset for action recognition," in *Proc. ICPR*, 2016.
- [20] A. Shahroudy, J. Liu, T.-T. Ng, and G. Wang, "Ntu rgb+ d: A large scale dataset for 3d human activity analysis," in *Proc. CVPR*, 2016, pp. 1010–1019.
- [21] M. E. Hussein, M. Torki, M. A. Gowayed, and M. El-Saban, "Human action recognition using a temporal hierarchy of covariance descriptors on 3d joint locations," in *Proc. IJCAI*, 2013, pp. 2466–2472.
- [22] M. A. Gowayed, M. Torki, M. E. Hussein, and M. El-Saban, "Histogram of oriented displacements (hod): describing trajectories of human joints for action recognition," in *Proc. IJCAI*, 2013, pp. 1351–1357.
- [23] F. Offi, R. Chaudhry, G. Kurillo, R. Vidal, and R. Bajcsy, "Sequence of the most informative joints (smij): A new representation for human skeletal action recognition," *J. Vis. Commun. Image Representation*, vol. 25, no. 1, pp. 24–38, 2014.
- [24] R. Chaudhry, F. Offi, G. Kurillo, R. Bajcsy, and R. Vidal, "Bio-inspired dynamic 3d discriminative skeletal features for human action recognition," in *Proc. CVPRW*, 2013, pp. 471–478.
- [25] M. Zanfir, M. Leordeanu, and C. Sminchisescu, "The moving pose: An efficient 3d kinematics descriptor for low-latency action recognition and detection," in *Proc. ICCV*, 2013, pp. 2752–2759.
- [26] J. Luo, W. Wang, and H. Qi, "Group sparsity and geometry constrained dictionary learning for action recognition from depth maps," in *Proc. ICCV*, 2013, pp. 1809–1816.
- [27] J. Wang, Z. Liu, Y. Wu, and J. Yuan, "Mining actionlet ensemble for action recognition with depth cameras," in *Proc. CVPR*, 2012, pp. 1290–1297.
- [28] M. Baccouche, F. Mamalet, C. Wolf, C. Garcia, and A. Baskurt, "Sequential deep learning for human action recognition," in *Human Behavior Understanding*, 2011, pp. 29–39.
- [29] U. Gaur, Y. Zhu, B. Song, and A. Roy-Chowdhury, "A 'string of feature graphs' model for recognition of complex activities in natural videos," in *Proc. ICCV*, 2011, pp. 2595–2602.
- [30] L. Wang and H. Sahbi, "Directed acyclic graph kernels for action recognition," in *Proc. ICCV*, 2013, pp. 3168–3175.
- [31] D. I. Shuman, S. K. Narang, P. Frossard, A. Ortega, and P. Vandergheynst, "The emerging field of signal processing on graphs: Extending high-dimensional data analysis to networks and other irregular domains," *IEEE Signal Process. Mag.*, vol. 30, no. 3, pp. 83–98, 2013.
- [32] M. Defferrard, X. Bresson, and P. Vandergheynst, "Convolutional neural networks on graphs with fast localized spectral filtering," in *Proc. NIPS*, 2016, pp. 3837–3845.
- [33] J. Bruna, W. Zaremba, A. Szlam, and Y. LeCun, "Spectral networks and locally connected networks on graphs," *arXiv preprint arXiv:1312.6203*, 2013.
- [34] M. Niepert, M. Ahmed, and K. Kutzkov, "Learning convolutional neural networks for graphs," in *Proc. ICML*, 2016.
- [35] D. K. Hammond, P. Vandergheynst, and R. Gribonval, "Wavelets on graphs via spectral graph theory," *Appl. Comput. Harmon. Anal.*, vol. 30, no. 2, pp. 129–150, 2011.
- [36] D. Bahdanau, K. Cho, and Y. Bengio, "Neural machine translation by jointly learning to align and translate," *arXiv preprint arXiv:1409.0473*, 2014.
- [37] J. K. Chorowski, D. Bahdanau, D. Serdyuk, K. Cho, and Y. Bengio, "Attention-based models for speech recognition," in *Proc. NIPS*, 2015, pp. 577–585.

- [38] K. Xu, J. Ba, R. Kiros, K. Cho, A. Courville, R. Salakhudinov, R. Zemel, and Y. Bengio, "Show, attend and tell: Neural image caption generation with visual attention," in *Proc. ICML*, 2015, pp. 2048–2057.
- [39] S. Hochreiter and J. Schmidhuber, "Long short-term memory," *Neural comput.*, vol. 9, no. 8, pp. 1735–1780, 1997.
- [40] A. Graves, "Generating sequences with recurrent neural networks," *arXiv preprint arXiv:1308.0850*, 2013.
- [41] N. Srivastava, E. Mansimov, and R. Salakhudinov, "Unsupervised learning of video representations using lstms," in *Proc. ICML*, 2015, pp. 843–852.
- [42] I. Sutskever, O. Vinyals, and Q. V. Le, "Sequence to sequence learning with neural networks," in *Proc. NIPS*, 2014, pp. 3104–3112.
- [43] J. Liu, A. Shahroudy, D. Xu, and G. Wang, "Spatio-temporal lstm with trust gates for 3d human action recognition," in *Proc. ECCV*, 2016, pp. 816–833.
- [44] W. Li, Z. Zhang, and Z. Liu, "Action recognition based on a bag of 3d points," in *Proc. CVPRW*, 2010, pp. 9–14.
- [45] P. Wang, W. Li, Z. Gao, C. Tang, J. Zhang, and P. Ogunbona, "Convnets-based action recognition from depth maps through virtual cameras and pseudocoloring," in *Proc. ACM MM*, 2015, pp. 1119–1122.
- [46] P. Wang, W. Li, Z. Gao, J. Zhang, C. Tang, and P. O. Ogunbona, "Action recognition from depth maps using deep convolutional neural networks," *IEEE Trans. Human Mach. Syst.*, vol. 46, no. 4, pp. 498–509, 2016.
- [47] O. Oreifej and Z. Liu, "Hon4d: Histogram of oriented 4d normals for activity recognition from depth sequences," in *Proc. CVPR*, 2013, pp. 716–723.
- [48] H. S. Koppula, R. Gupta, and A. Saxena, "Learning human activities and object affordances from rgb-d videos," *Int. J. Robot. Res.*, vol. 32, no. 8, pp. 951–970, 2013.
- [49] J. Sung, C. Ponce, B. Selman, and A. Saxena, "Unstructured human activity detection from rgb-d images," in *Proc. ICRA*, 2012, pp. 842–849.
- [50] V. Bloom, V. Argyriou, and D. Makris, "Dynamic feature selection for online action recognition," in *Proc. HBU*, 2013, pp. 64–76.
- [51] V. Bloom, D. Makris, and V. Argyriou, "G3d: A gaming action dataset and real time action recognition evaluation framework," in *Proc. CVPRW*, 2012, pp. 7–12.
- [52] B. Ni, G. Wang, and P. Moulin, "Rgb-d-hudaact: A color-depth video database for human daily activity recognition," in *Proc. ICCVW*, 2011, pp. 1147–1153.
- [53] C. Chen, R. Jafari, and N. Kehtarnavaz, "Utd-mhad: A multimodal dataset for human action recognition utilizing a depth camera and a wearable inertial sensor," in *Proc. ICIP*, 2015, pp. 168–172.
- [54] L. Wang, J. Zhang, L. Zhou, C. Tang, and W. Li, "Beyond covariance: Feature representation with nonlinear kernel matrices," in *Proc. ICCV*, 2015, pp. 4570–4578.
- [55] Z. Huang and L. Van Gool, "A riemannian network for spd matrix learning," in *Proc. AAAI*, 2017.
- [56] M. T. Harandi, M. Salzmann, and R. Hartley, "From manifold to manifold: Geometry-aware dimensionality reduction for spd matrices," in *Proc. ECCV*, 2014, pp. 17–32.
- [57] O. Tuzel, F. Porikli, and P. Meer, "Region covariance: A fast descriptor for detection and classification," in *Proc. ECCV*, 2006, pp. 589–600.
- [58] Z. Huang, C. Wan, T. Probst, and L. Van Gool, "Deep learning on lie groups for skeleton-based action recognition," *arXiv preprint arXiv:1612.05877*, 2016.
- [59] Y. Yang, C. Deng, S. Gao, W. Liu, D. Tao, and X. Gao, "Discriminative multi-instance multitask learning for 3d action recognition," *IEEE Trans. Multimedia*, vol. 19, no. 3, pp. 519–529, 2017.
- [60] J.-F. Hu, W.-S. Zheng, J. Lai, and J. Zhang, "Jointly learning heterogeneous features for rgb-d activity recognition," in *Proc. CVPR*, 2015, pp. 5344–5352.
- [61] G. Evangelidis, G. Singh, and R. Horaud, "Skeletal quads: Human action recognition using joint quadruples," in *Proc. ICPR*, 2014, pp. 4513–4518.
- [62] T. Batabyal, T. Chattopadhyay, and D. P. Mukherjee, "Action recognition using joint coordinates of 3d skeleton data," in *Proc. ICIP*, 2015, pp. 4107–4111.
- [63] M. Ranzato and G. E. Hinton, "Modeling pixel means and covariances using factorized third-order boltzmann machines," in *Proc. CVPR*, 2010, pp. 2551–2558.