

HMANet: Hyperbolic Manifold Aware Network for Skeleton-Based Action Recognition

Jinghong Chen^{1b}, Chong Zhao^{1b}, Qicong Wang^{1b}, and Hongying Meng^{1b}, *Senior Member, IEEE*

Abstract—Skeleton-based action recognition has attracted significant attentions in recent years. To model the skeleton data, most popular methods utilize graph convolutional networks to fuse nodes located in different parts of the graph to obtain rich geometric information. However, these methods cannot be generalized to different graph structures due to their dependencies on the input of the topological structure. In this article, we design a novel hyperbolic manifold aware network without introducing a dynamic graph. Instead, it leverages Riemannian geometry attributes of a hyperbolic manifold. Specifically, this method utilizes the Poincaré model to embed the tree-like structure of the skeleton into a hyperbolic space to automatically capture hierarchical features, which may explore the underlying manifold of the data. To extract spatio-temporal features in the network, the features in manifold space are projected to a tangent space, and a tangent space features translation method based on the Levi-Civita connection was proposed. In addition, we introduce the geometric knowledge of Riemannian manifolds to further explain how features are transformed in the tangent space. Finally, we conduct experiments on several 3-D skeleton data sets with different structures, successfully verifying the effectiveness and advancement of the proposed method.

Index Terms—Action recognition, hyperbolic manifold, Poincaré model, Riemannian geometry, spatio-temporal features.

I. INTRODUCTION

ACTION recognition is one of the most important fields in computer vision research. It utilizes computer vision methods to determine the action category of the camera recorded data. Generally, two different types of data are used in action recognition tasks, RGB video data, and skeleton data. Most methods based on RGB video data use convolutional neural networks (CNNs) to extract image information or use traditional methods to extract video optical flow trajectory information. RGB video data has the advantages of easy

collection and data regularization, but it is susceptible to interference from the shooting environment. With the development of reliable skeleton estimation methods in depth video [1] and RGB video [2], the 3-D joint positions of human bones in action videos can be easily obtained in real time, which greatly promotes the research and application of skeleton-based action recognition. In recent years, skeleton-based human action recognition has received widespread attention. It mainly uses the Euclidean coordinates of 3-D joint points for modeling. The compact skeleton data makes the model more efficient and robust to changes in perspective and environment. These methods have achieved desirable results.

The methods based on manual features [3]–[6] capture spatio-temporal or geometric features of the skeleton sequence, while the methods based on deep learning is directly supervised by the action category to learn discriminative spatio-temporal features. Although action recognition methods based on manual features can usually achieve good performance, these methods have intrinsic limitations, especially that they can only extract shallow features. Deep learning provides a way to obtain high semantic representations. For example, taking advantage of the characteristics of RNN being suitable for time series data processing, methods based on RNN have been proposed to improve the ability to learn temporal context. Thus, long short-term memory (LSTM) was introduced to extract time-series features. Zhang *et al.* [7] applied geometric joint features to multilayer LSTM networks instead of joint positions. Ma *et al.* [8] utilized dynamic evolution of time series by introducing differences of time series as inputs to the LSTM. The main drawback of these methods is that they lack spatial modeling capabilities, resulting in poor results. CNN has an excellent ability to extract high-level semantic information. Many approaches [9]–[11] have utilized the CNN model for action recognition by encoding skeletal joints as pseudo images, and then input it into the network. Zhang *et al.* [11] mapped a skeleton sequence to an image to facilitate spatio-temporal modeling by CNN. Banerjee *et al.* [12] proposed four feature representations of the sequence of key joints, and utilized CNNs to encode these features for classification. Compared to RNN, a significant challenge of using CNN is how to organize sequential data for the natural input of the model. Most methods directly convert the skeleton data into images, which may lead to the loss of spatial information and usually complex calculations, limiting their practical applications. Therefore, the application of CNN for skeleton-based action recognition is still an

Manuscript received 25 November 2021; revised 10 March 2022; accepted 26 April 2022. Date of publication 2 May 2022; date of current version 12 June 2023. This work was supported by the Shenzhen Science and Technology Program under Grant JCYJ20200109143035495. (Corresponding author: Qicong Wang.)

Jinghong Chen and Qicong Wang are with the Shenzhen Research Institute, Xiamen University, Shenzhen 518057, China, and also with the Department of Computer Science and Technology, Xiamen University, Xiamen 361005, China (e-mail: qcwang@xmu.edu.cn).

Chong Zhao is with the Department of Computer Science and Technology and the Xiamen Deep Geometry Intelligent Technology Research Institute, Xiamen University, Xiamen 361005, China.

Hongying Meng is with the Department of Electronic and Computer Engineering, Brunel University, London UB8 3PH, U.K. (e-mail: hongying.meng@brunel.ac.uk).

Color versions of one or more figures in this article are available at <https://doi.org/10.1109/TCDS.2022.3171550>.

Digital Object Identifier 10.1109/TCDS.2022.3171550

unsolved research problem. In recent years, graph-based skeleton action recognition has become a research hotspot in the field of computer vision due to its excellent performance. The ST-GCN proposed in [13] applied graph convolutional networks (GCNs) to action recognition for the first time which achieved a great improvement in accuracy. In recent researches, Li *et al.* [14] modeled the skeleton sequence as a graph, and applied GCN to capture spatial and temporal dynamics to provide high performance. Although GCN-based methods have achieved excellent accuracy, they have limited applications and are costly in terms of memory when there are larger numbers of nodes. In addition, they cannot be generalized to different structure of graph.

All the previous approaches are defined in the Euclidean space. However, the underlying anatomical structure of the data often contains more geometric information in non-Euclidean spaces, so Euclidean space may not be the best choice for modeling hierarchical data. Recent studies have proved that complex types of data (such as graph data) in many fields exhibit topological structures that are closely related to manifolds. Under such circumstances, the Euclidean space cannot provide maximum expression ability or meaningful geometric representation. For example, Sala *et al.* [15] proved that arbitrary tree structures cannot be embedded with arbitrary low distortion (i.e., almost preserving their metric) in the Euclidean space with infinite dimensions, but this task becomes strikingly easy in the hyperbolic space with only two dimensions where the exponential growth of distances matches the exponential growth of nodes with the tree depth. Therefore, neural network operations defined directly in the data-related space [16] may benefit the learning process. Different from learning joints embedding directly in the Euclidean space, we explore the modeling space of the skeleton graph sequence in the non-Euclidean geometry. However, deep learning in these non-Euclidean spaces has been rather limited, the main reason being the nontrivial or impossible principled generalizations of basic operations (e.g., vector addition, matrix-vector multiplication, vector translation, and vector inner product). Thus, classic tools, such as feedforward networks or recurrent networks have no corresponding representations in these spaces, and it is difficult to find natural mathematical descriptions for basic operations such as convolution. Inspired by research [17], we get an idea from the bijection between the hyperbolic space and the tangent space. The classic operations can be generalized to tangent spaces through the logarithmic map. In this way, the spatio-temporal features can be obtained by applying Euclidean filters on feature map in tangent space. In this article, we construct a 3-D action recognition framework (HMANet) that leverages hyperbolic space to make spatio-temporal features full of hierarchy. Our contributions can be summarized as follows.

- 1) To the best of our knowledge, our HMANet introduces a hyperbolic manifold into the field of 3-D action recognition for the first time. It devotes to mining the spatial configuration of the skeleton sequence. For features represented in the hyperbolic space, we mix temporal and spatial filters to extract spatio-temporal features in the tangent space.

- 2) Explain how our network learns the features in the tangent space from the perspective of differential geometry, and establish relationship between the metric tensor of the Riemannian manifold and the features in the tangent space through mathematical theory. The introduction of manifold theory into the model makes it more explanatory.
- 3) Propose a hyperbolic aware bias for features in the tangent space of manifold. It utilizes the parallel transport with respect to the Levi-Civita connection to translate the tangent vector along the geodesic to make the captured features lie in different tangent spaces of manifold, such that the model can automatically aware of underlying manifold.

The remainder of this article is organized as follows. Section II reviews the related approaches and discusses their relationships to the present works. Section III gives a detailed description of our method and the corresponding network architecture, while supplying a theoretical analysis. Comprehensive experimental results and analysis are provided in Section IV, and finally, a conclusion is drawn in Section V.

II. RELATED WORK

A. Skeleton-Based Action Recognition With CNN

Most of the methods based on CNN flatten the 3-D skeleton sequence into pseudo images with joints and frames as different dimensions, and the feature learning follows the methods in image. Li *et al.* [18] encoded the pairwise distances between joints into RGB images, and separately trained CNN models in four orthogonal planes with empirical fusion schemes account for view invariance. Banerjee *et al.* [12] proposed a CNN model, which leverages features estimated from angular information and kinematics of human to capture complementary characteristics of the sequence of key joints. The approach mentioned in [19] is a CNN-based method that utilizes a gating mechanism for images generated from a specific order of skeletons. The two-stream attention mask in CNN was reported in [20]. Li *et al.* [21] used the features in methods [9], [18] and the LSTM network to study the multiclassifier classification model of the maximum, multiplicative and average decision score fusion scheme. These methods are not sensitive to subtle movement changes within the class which can be rectified by using more specialized features. Huynh-The *et al.* [22] studied specialized geometric feature extraction techniques, including joint orientation, which provided impressive performance. Recent methods [23], [24] exploit transition geometric features alongside frame-wise geometrical features, which is a very crucial step toward utilizing motion information. The features learned by these methods treat the data as an image, and thus fail to effectively express the long-distance interaction relationship in the skeleton. Although the CNN operator can indeed form an overall feature representation through the local convolution kernel, it neglects the interaction of the long-distance joints. Moreover, the Euclidean distance between joint coordinates cannot accurately describe their geometric distance. For the purpose of learning the features implying underlying manifold,

our method attempts to mine this geometric topology in hyperbolic space, which enables the distance between coordinates to express their geometric structure to a certain extent.

B. Representations in Non-Euclidean Space

In order to explore more robust skeleton features in non-Euclidean space, one approach is to directly employ manifold data as the original input. For example, researchers express rotation relationships as points in the Lie group $SO(3)$, and describe the skeleton motion information through the rotation relationships between each pair of 3-D vectors, so as to eliminate the influence of viewing angle changes and learn more robust features. Vemulapalli *et al.* [25] first proposed performing action recognition by using $SO(3)$ to represent human bones (rotation and translation), LieNet [26] further realized deep learning curve clusters by defining rotation map transformation, Vemulapalli and Chellapa [27] introduced the concept of a rolling map in mathematics, which mapped the $SO(3)$ representation of the human skeleton to the tangent space, and utilized SVM for linear classification. These methods manually characterize the data in a specific manifold, which may lose part of the original information, resulting in poor results, and networks specially designed for them often bring a large amount of calculation.

Another more reasonable approach is to generalize the deep neural network to the non-Euclidean geometry. Specifically, it uses deep learning to automatically embed the data on the Riemannian manifold. For example, in order to construct a model on a Riemannian manifold, Mathieu *et al.* [28] proposed Poincaré variational autoencoder and showed a better generalization for hierarchical structures. In this article, we focus on hyperbolic manifolds, which is a non-Euclidean space with a constant negative Gaussian curvature and has the ability to efficiently model hierarchical structures. In machine learning, hyperbolic representations greatly outperformed Euclidean embeddings for hierarchical, taxonomic, or entailment data recently. Disjoint subtrees from the implicit hierarchical structure are well clustered in the embedding space. However, appropriate deep learning tools are needed to embed feature data in this space and use it in downstream tasks. Ganea *et al.* [29] established the connection between the hyperbolic manifold and Euclidean space in the context of neural network and deep learning, and generalizes basic operators, polynomial regression, and feedforward network to the Poincaré model of a hyperbolic manifold. Ungar [30] combined the gyrovector space and the generalized Möbius transformation with the popular properties of the Riemannian geometric, smoothly parametrize basic operations and objects in all spaces of constant negative curvature using a unified framework that depends only on the curvature value. Then, the Euclidean space and hyperbolic spaces can be continuously deformed into each other.

C. Neural Networks on Hyperbolic Manifold

Recently, there have been some attempts to design neural networks in hyperbolic space. Specifically, the pioneering research on learning representation in hyperbolic spaces was

reported in [31]. Then, in the research [29], hyperbolic neural networks were introduced, linking hyperbolic geometry with deep learning. Subsequent related works provided analogies on the hyperbolic manifolds of classic operations, or developed several other algorithms, such as Poincaré GloVe [32] and hyperbolic aware mechanism networks [17]. In addition, their method is also more general for graph sequence data because they are naturally in the non-Euclidean space. Chami *et al.* [33] utilized hyperbolic geometry to construct a graph neural network. Considering that there is a bijection between the hyperbolic space and the tangent space, scholars can first perform the convolution operation on the tangent space, and then project the extracted features back as a trajectory on the manifold. Since the hyperbolic distance between unrelated samples in a hyperbolic manifold will increase exponentially than the distance between similar samples, it may be better to construct a classification model for a human skeleton on a hyperbolic manifold. To this end, we are dedicated to propose a spatio-temporal manifold-aware network for a specific model of hyperbolic geometry (i.e., the Poincaré model). This network does not generate node embeddings by inputting human spatio-temporal graph, but explores more reasonable manifold projections, such that the projection features are more discriminative and the network can be generalized to different skeleton structures. In addition, regarding the interpretability of neural networks, Hauser *et al.* [34] took feature transformation as the transformation of Riemannian metric tensor on manifold from the perspective of differential geometry. This article attempts to study these issues on hyperbolic manifolds, and combines the network architecture with hyperbolic space, taking advantages of its good hierarchical structure modeling capabilities to further strengthen the exploration of hierarchical structures.

III. PROPOSED METHOD

In this section, we describe our classification model HMANet in detail. The framework is shown in Fig. 1. The convolutional layer of our network consists of a spatial filter and a temporal filter. Such blocks are added to capture the spatio-temporal features of the skeleton sequence. The model first expands the joint dimension by affine transformation, leveraging exponential function to map coordinates to hyperbolic space and renew position coordinates, then performs affine transformation in the Euclidean space by the logarithmic map. We will describe important components of our framework in the following sections in details. Table I shows the used notations and their corresponding definitions.

A. Poincaré Model of Hyperbolic Geometry

3-D human skeleton can be represented as a graph composed of nodes and edges due to the spatial topology of joints. Traditional Euclidean space is a linear manifold, any parameterization method does not have the ability to represent a graph. Nevertheless, hyperbolic space provides a more reasonable embedding for human joints due to its unique metric properties. Cannon *et al.* [35] gave five isometric models of hyperbolic space. To represent hyperbolic space in a simple

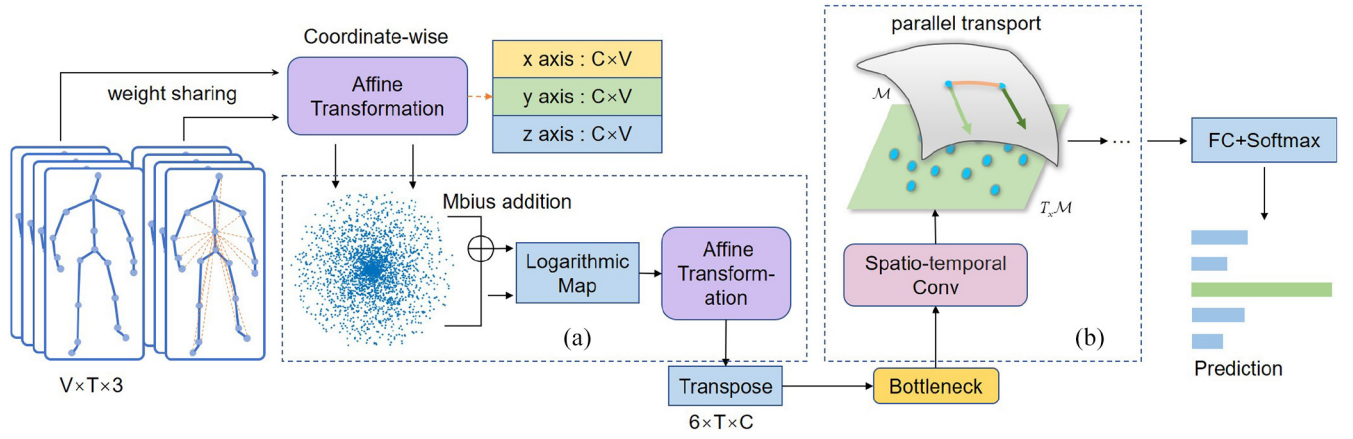


Fig. 1. Illustration of our framework (HMANet). There are mainly two stages in our framework, including (a) coordinate-wise affine transformation and (b) convolution equipped with manifold transaction. At the first stage, we concatenate position tensor with the corresponding deviation tensor and utilize affine transformation to transform coordinates of joints for each dimension. We stack several layers, at the start of next layer, we map the two to hyperbolic space and perform Möbius addition, then utilize the logarithmic map to map them back to Euclidean space. The bottleneck followed by is to point-wise expand the dimension. In stage (b), we adopt manifold transaction in convolution layer to make it manifold aware.

TABLE I
NOTATIONS AND DEFINITIONS

Notations	Definitions
\mathcal{M}	a smooth manifold
H_2	a 2D Poincaré disc
δ_H	the metric in hyperbolic space
δ_E	the metric in Euclidean space
\mathcal{D}^n	an n-dimensional open unit ball in Euclidean space
γ_x	the Riemannian metric tensor at point x of manifold
I_n	the n-order identity matrix
λ_x	the conformal factor on manifold
\oplus	the Möbius addition on Poincaré model
V	the number of joint points
T	the number of frames in one action
J_t^v	the 3D coordinate of the v -th joint in the t -th frame
w_t^c	the position vector of the c -th channel in the t -th frame
\tilde{w}_t^c	the deviation vector of the c -th channel in the t -th frame
$T_x \mathcal{D}^n$	the tangent space at point x on manifold \mathcal{D}^n
U_x	an open set on manifold \mathcal{M}
φ_x	A coordinate function that maps elements in U_x to \mathcal{R}^n
E_x	a set of basis vectors in tangent space
$T_x^* \mathcal{M}$	the cotangent space of the manifold \mathcal{M}
D_v	the directional derivative of the direction v
E_x^*	a set of basis vectors in cotangent space
v^l	a feature vector of layer l
\mathcal{H}^l	the Jacobian matrix of mapping between two manifolds
\mathcal{J}	a smooth second order tensor on manifold

way, we choose to study in the Poincaré model. In Fig. 2, any two gesture joints can be embedded as two points x^1, x^2 in the Poincaré disc H_2 . Poincaré disc $H_2 := \{(x_1, x_2) \mid x_1^2 + x_2^2 < 1\}$ is a 2-D case of hyperbolic geometry, wherein the distance metric changes. Near each joint, the metric is related to the position of the node, whereas the shortest path between two nodes is not straight-line distances. We show that this distance can reflect the topological structure of the joints with the help of the following definitions and formulas.

Hyperbolic space is a Riemannian manifold with constant negative curvature, which is a curved metric space that does not have a distance-preserving relationship with the Euclidean space. For nodes represented in hyperbolic space, the metric space expands exponentially with the distance from the original point. Thus it is more advantageous to represent hierarchical information. The Poincaré sphere model (\mathcal{D}^n, γ) is an

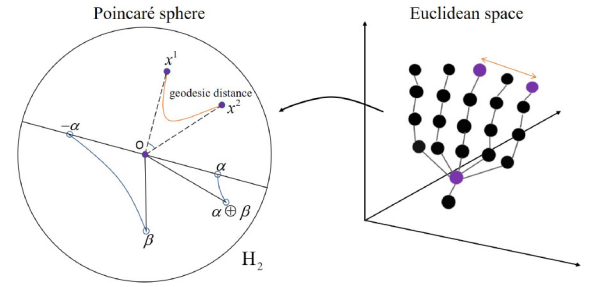


Fig. 2. Left is the 2-D Poincaré disk embedded in the Euclidean space, where the distance between every two points is the geodesic distance. The right is an illustration of the hand gesture. Embed any two joints into the disk, the distance between them is converted into hyperbolic distance.

n -dimensional hyperbolic space equipped with a Riemannian metric γ , defined as $\mathcal{D}^n = \{x \in \mathcal{R}^n : \|x\| < 1\}$. The Riemannian metric $\gamma_x : T_x \mathcal{D}^n \times T_x \mathcal{D}^n \rightarrow \mathcal{R}$ is a family of positive-definite quadratic forms that smoothly vary with point x on the manifold

$$\gamma_x = \lambda_x^2 \cdot \gamma^E \quad (1)$$

where $\lambda_x = [2/(1 - \|x\|^2)]$ is the conformal factor and $\gamma^E = I_n$ is the Euclidean metric tensor. Therefore, the metric of 2-D disc space is defined as

$$ds^2 = \left(\frac{2}{1 - x_1^2 - x_2^2} \right)^2 (dx_1^2 + dx_2^2). \quad (2)$$

It can be seen from the formula that the closer the point is to the edge of the disc, the greater the distance represented by the coordinate difference $(\Delta x_1, \Delta x_2)$. Take the gesture skeleton as an example, we illustrate the inspiration for the key technology of embedding skeleton joints in hyperbolic space. In Fig. 2, we use a geodesic to show the hyperbolic distance between two points. For the joint points x^1, x^2 in the disc, the geodesic distance is defined as follows:

$$\delta_H(x^1, x^2) = \cosh^{-1} \left(1 + 2 \frac{\delta_E(x^1, x^2)^2}{(1 - \|x^1\|^2)(1 - \|x^2\|^2)} \right) \quad (3)$$

where δ_H represents the hyperbolic distance and δ_E represents the Euclidean distance, they can be extended to the case of the 3-D skeleton data. Suppose $\|x^1\| = \|x^2\| = \tau$, there is

$$\lim_{\tau \rightarrow 1} \delta_H(x^1, x^2) = \delta_H(x^1, 0) + \delta_H(x^2, 0). \quad (4)$$

In other words, the shortest path between x^1 and x^2 is almost the same as the path through the origin. This is analogous to a tree structure, in which the shortest path between two sibling nodes is the path through their parent node. The tree-like property of hyperbolic space is a key attribute for feature embedding. Given any two points on the disc, no matter how small the angle between them to the center is, this property can be satisfied. Therefore, the hyperbolic distance can well reflect the distance in the sense of joint topology, the natural embedding of the hierarchical structure can be found in the hyperbolic space.

Hyperbolic space is a nonlinear space, thus the addition defined in hyperbolic space is different from the Euclidean space. It is called Möbius addition, denoted as \oplus . For any two points α and β in disc

$$\alpha \oplus \beta = \frac{(1 + 2\langle \alpha, \beta \rangle + \|\beta\|^2)\alpha + (1 - \|\alpha\|^2)\beta}{1 + 2\langle \alpha, \beta \rangle + \|\alpha\|^2\|\beta\|^2} \quad (5)$$

Fig. 2 describes the operation from a geometric perspective. As shown in the figure, $\alpha \oplus \beta$ is obtained by translating the triangle along the side $-O\alpha$. Connection can be established by two congruent triangles

$$\begin{cases} d(-\alpha, \beta) = d(0, \alpha \oplus \beta) \\ d(0, \beta) = d(\alpha, \alpha \oplus \beta). \end{cases} \quad (6)$$

Therefore, combined with the hyperbolic distance formula, it can be known that when β is closer to the center, $\alpha \oplus \beta$ is closer to α , and when β is closer to the edge, the coefficient in the α direction is greater. Simultaneously, if the directions of α and β are closer, $\alpha \oplus \beta$ is farther from the center.

B. Architecture of HMANet

For the input motion coordinates, we propose an end-to-end deep learning framework. We first represent the skeleton sequence with V joints and T frames as a tensor of shape $V \times T \times 3$. For the skeleton of a person in frame t , we formulate it as $J_t = (J_t^1, J_t^2, \dots, J_t^V)^T$ and $J_t^v = (J_{tx}^v, J_{ty}^v, J_{tz}^v)$ is the 3-D joint coordinates. In this way, the skeleton sequence is regarded as an image with three channels. Considering that: 1) The two dimensions of the image represent joints and frames, which are usually not equivalent and 2) the movement of a joint is not only related to the local area, but also related to the distant joints. We treat each joint of the skeleton as a channel, and learn the global response of all channels through affine transformation. However, any two channels of the output feature are no longer in a parallel relationship, and they share part of the same information. To this end, we propose a method of transforming features through hyperbolic space, such that the distance of features in new space more accurately reflects their relevance.

Suppose that the coordinate of the human center of gravity in the skeleton is J_t^1 , the difference of the coordinates

$\tilde{J}_t = J_t - J_t^1$ is calculated in each frame, and these 3-D vectors form a tensor with the shape of $V \times T \times 3$, which is called the deviation tensor. We divide the skeleton sequence into three parts according to the coordinate dimension, and connect each part with the deviation tensor according to the corresponding dimension to obtain three tensors with the shape of $V \times T \times 2$. After that, we use three affine transformations to independently aggregate the global features of all joints for each dimension of the coordinate, and use the batchnorm to normalize them, then connect the three dimensions. The obtained tensor is composed of position vectors and deviation vectors. Repeating the learning of the global features, before each subsequent affine transformation, we use the following method to renew the position vectors.

Let $W \in R^{C \times T \times 6}$ be the output of the first layer, where C represents the number of channels. The component of the output tensor at frame t is written as $W_t = (w_{tx}, w_{ty}, w_{tz}, \tilde{w}_{tx}, \tilde{w}_{ty}, \tilde{w}_{tz}) \in R^{C \times 6}$, and channel c contains a position vector $w_t^c = (w_{tx}^c, w_{ty}^c, w_{tz}^c)$, and a vector $\tilde{w}_t^c = (\tilde{w}_{tx}^c, \tilde{w}_{ty}^c, \tilde{w}_{tz}^c)$ obtained by affine transformation of the original tensor. Let $A_j \in R^{C \times V}$ and $b_j \in R^C$ ($j = x, y, z$) be optimizable parameters, the output of each coordinate dimension is obtained by affine transformation

$$\begin{cases} w_{tx} = A_x J_{tx} + b_x, & \tilde{w}_{tx} = A_x \tilde{J}_{tx} + b_x \\ w_{ty} = A_y J_{ty} + b_y, & \tilde{w}_{ty} = A_y \tilde{J}_{ty} + b_y \\ w_{tz} = A_z J_{tz} + b_z, & \tilde{w}_{tz} = A_z \tilde{J}_{tz} + b_z \end{cases} \quad (7)$$

In order to map them to points in hyperbolic space, we refer to the bijection of hyperbolic space and tangent space at one point proposed by Ganea *et al.* [29], called exponential map and logarithmic map. The following are the definitions of exponential map and logarithmic map. $\forall x \in \mathcal{D}^n$

$$\exp_x(v) = x \oplus \left(\tanh\left(\frac{\lambda_x \|v\|}{2}\right) \phi(v) \right) \quad (8)$$

$$\log_x(y) = \frac{2}{\lambda_x} \tanh^{-1}(\| -x \oplus y \|) \phi(-x \oplus y) \quad (9)$$

where $\phi(r) = (r/\|r\|)$ represents vector unitization. We pay attention to case $x = 0$, use the projection function to map the position vectors and the deviation vectors to the hyperbolic space, and perform the Möbius addition to renew the position vectors. However, hyperbolic space is a nonlinear space, and affine transformations cannot be directly applied to the features in it. To this end, we convert the features to the tangent space. Tangent space is the linearization of the manifold, which can be regarded as the set of all differential operators at a point on manifold. Since it is a vector space, and there is a bijection between it and manifold, the manifold features can be projected to the tangent space of a point to preserve the manifold structure. We use the logarithmic map to project the vectors from the manifold to the tangent space, such that the loss function in the Euclidean space can be employed to optimize the model

$$w_t^c \leftarrow \log_0(\exp_0(w_t^c - \tilde{w}_t^c) \oplus \exp_0(\tilde{w}_t^c)). \quad (10)$$

We use the deviation vectors to renew the position vectors in the hyperbolic space. Based on the previous discussion, if the included angle of any two deviation vectors is small,

the newly obtained corresponding position vector angle will become smaller. Besides, the larger the norm of the deviation vector, it means that in the corresponding position vector, the larger weights are more likely derived from the neighboring points. Therefore, the feature independence is stronger, and the coefficients of this direction are also larger. This is intuitive in the feature space.

Finally, we obtain a tensor with a shape of $C \times T \times 6$. We designate the dimensions of the tangent vectors as channels by transpose, and use bottleneck to increase the feature dimensions before sending it to the convolution layer, which is equivalent to obtaining the coordinate representation of the high-dimensional manifold through manifold immersion. Then we combine the spatial and temporal filters to extract high-order features. Fig. 1 illustrates the entire network framework.

C. Convolution Layer on Tangent Space of Manifold

To discuss the deep convolution block in our model HMANet, we introduce some knowledge about tangent space in this section. As mentioned in Section III-A, the Riemannian metric is a quadratic form acting on the tangent space, which can induce the geodesic distance on the manifold.

Metric \tilde{g} and metric g are conformal when they define the identical angle. The Poincaré sphere and Euclidean space are conformal, namely $\forall x \in \mathcal{D}^n, u, v \in T_x \mathcal{D}^n \setminus \{0\}$, there is

$$\cos(\angle(u, v)) = \frac{g_x^D(u, v)}{\sqrt{g_x^D(u, u)}\sqrt{g_x^D(v, v)}} = \frac{\langle u, v \rangle}{\|u\|\|v\|}. \quad (11)$$

Therefore, the length of the tangent vector can be naturally defined as the length in Euclidean space.

In order to transform the manifold features in deep learning based on affine transformation, we map the learned manifold features to the tangent space, thereby obtaining a tangent vector field on the manifold. Since Riemannian manifold \mathcal{M} has a local Euclidean structure, there is a family of coordinate charts $\{(U_x, \varphi_x)\}$ that form an open cover of \mathcal{M} , and each coordinate function $\varphi_x \in C_x^\infty$ is a homeomorphism from U_x to an open set of \mathcal{R}^n . Given a point $x \in \mathcal{M}$ which is mapped to \mathcal{R}^n by the coordinate function φ_x , its tangent space $T_x \mathcal{M}$ has a set of natural basis $E_x = \{(\partial/\partial x_1), \dots, (\partial/\partial x_n)\}$. $\forall v \in T_x \mathcal{M}$, the derivative of the function along the direction v at x is: $D_v[\varphi] = v \cdot d\varphi$, where $[\varphi]$ represents a germ of function at x , that is, all equal functions in a sufficiently small neighborhood. We call $d\varphi$ the cotangent vector, the space $T_x^* \mathcal{M}$ composed of cotangent vectors is called the dual space.

The neural network applies linear transformation to the functional operator in the tangent space through affine transformation. Specifically, considering the cotangent space $T_x^* \mathcal{M}$, there is a dual natural basis $E_x^* = \{dx^1, \dots, dx^n\}$. We give a smooth tensor $\mathcal{J} : T_x \mathcal{M} \times T_x^* \mathcal{M} \rightarrow \mathcal{R}$ of type $(1, 1)$ which contains a family of Riemannian metrics on \mathcal{M} , they can perform inner product on the tangent vector as a special quadratic form. In addition, the Riemannian metric tensor is an endomorphism of the tangent bundle (i.e., $\mathcal{J} : T_x \mathcal{M} \rightarrow T_x \mathcal{M}$). The weight matrix of the neural network implies the metric tensor \mathcal{J} . Generalizing to a more general case, when the feature dimension is expanded by the network layer, the underlying

manifold of features is immersed into the higher-dimensional manifold.

D. Hyperbolic Aware Bias Based on Levi-Civita Connection

Since the tangent space at each point of manifold is not identical, while the captured features are projected to the tangent space of the original point (i.e., $T_0 \mathcal{M}$) through the logarithmic map, we introduce the following theorem and propose a bias in tangent space to transfer the translation along the geodesic of manifold to the tangent space, converting the tangent vector $v \in T_0 \mathcal{M}$ to a tangent vector $v' \in T_x \mathcal{M}, x \neq 0$.

As referred in [29], in the manifold (\mathcal{D}^n, g) , the parallel transport w. r. t. the Levi-Civita connection of a vector $v \in T_0 \mathcal{D}^n$ to another tangent space $T_x \mathcal{D}^n$ is given by the following isometry:

$$P_{0 \rightarrow x}(v) = \log_x(x \oplus \exp_0(v)) = \frac{\lambda_0}{\lambda_x} v = (1 - \|x\|^2)v. \quad (12)$$

This equation is crucial for defining and optimizing the parameters shared between different tangent spaces. Back to our network, the features are regarded as the vectors in the tangent space of original point $T_0 \mathcal{M}$ through the logarithmic map. A bias is applied to each feature vector to control the distance from the origin point. The distance defined according to the tangent space may reflect its hyperbolic nature. For this reason, we employ the above function to translate the tangent space features. Given a feature in vector form $v_{ti} \in T_x \mathcal{D}^n$, where t represents the temporal dimension and i represents the spatial dimension, we set an optimizable bias parameter $b \in T_0 \mathcal{D}^n$ to translate the features in tangent space. Combining (8) and (12), the conformal factor of the feature v_{ti} mapping onto the manifold can be obtained by the tanh function. Therefore, the deformation of bias b on each feature is

$$P_{0 \rightarrow \exp_0(v_{ti})}(b) = (1 - (\tanh\|\exp_0(v_{ti})\|)^2) \cdot b. \quad (13)$$

Then, we write the convolution operation on the l th layer in the network into the following form:

$$v_{ii}^{l+1} = \sigma \left(P_{0 \rightarrow \exp_0(v_{ii}^l)}(b) + (\mathcal{F}_t(v_{ii}^l) + \mathcal{F}_i(v_{ii}^l)) \right) \quad (14)$$

where \mathcal{F}_t and \mathcal{F}_i represent temporal convolution and spatial convolution, respectively. After adding the deformed bias, we make features adapt to different tangent spaces, such that the captured features are located in the tangent space at different points of the manifold, and the features can be represented in the manifold space more accurately. We will demonstrate the effectiveness of this operation by experiments in Section IV.

E. Submanifold Immersion and Feature Embedding in Hyperbolic Space

In this section, we explain that after the dimensionality of the features in tangent space is increased by bottleneck, the obtained features can be regarded as vectors in the tangent space of a high-dimensional manifold. In deep learning based on the Euclidean geometry, the features of the neural network can be regarded as a set of Cartesian coordinates in the Euclidean space. The tangent space is a local linear approximation of the manifold, which can be parameterized by the

coordinates in the Euclidean space. For the features in manifold, we project them to the tangent space by the logarithmic map. In this way, we can learn the tangent space features by performing a linear transformation on the coordinates. If the dimension of features in each layer is constant, the coordinates can be renewed by a full-rank Jacobian matrix, and a positive-definite Riemannian metric is maintained. However, the dimension of features in the actual network increases as the layers deepens. From the perspective of differential geometry, if the rank of the map Jacobian matrix is equal to the dimension before the map, the manifold can be immersed in a higher-dimensional space.

In our network, in order to more easily classify the data, we embed features into higher-dimensional hyperbolic manifolds. Let \mathcal{M} and \mathcal{N} be m -dimensional and n -dimensional smooth manifolds, respectively, where $m \leq n$: $\mathcal{M} \rightarrow \mathcal{N}$ is a smooth map on manifold. Each tangent space feature in the layer l is represented as $v^l = \log_0(x^l) \in T_0\mathcal{D}^m$ by the logarithmic map on the manifold \mathcal{M} , and $v^{l+1} = \log_0(x^{l+1}) \in T_0\mathcal{D}^n$ on the manifold \mathcal{N} . Denoting that

$$\begin{aligned} h^l(v^l) &:= (\log_0 \circ f \circ \log_0^{-1})(v^l) \\ &= (\log_0 \circ f \circ \exp_0)(v^l) \end{aligned} \quad (15)$$

where $v^l \in \mathcal{R}^m$, $h^l(v^l) \in \mathcal{R}^n$, the Jacobian matrix of map f is given by two coordinate functions

$$\text{Jacobi}(h^l) = \begin{bmatrix} \frac{\partial h_1^l}{\partial v_1^l} & \cdots & \frac{\partial h_1^l}{\partial v_m^l} \\ \vdots & \ddots & \vdots \\ \frac{\partial h_n^l}{\partial v_1^l} & \cdots & \frac{\partial h_n^l}{\partial v_m^l} \end{bmatrix}. \quad (16)$$

Denoting this matrix as \mathcal{H}^l , the network will learn this matrix to ensure

$$\text{rank}(h^l) := \text{rank}(\mathcal{H}^l) = m. \quad (17)$$

In this case $\forall x \in \mathcal{M}$, the Jacobian matrix of the coordinate function h^l is general nondegenerate, then f is an immersion of smooth manifold \mathcal{M} in \mathcal{N} . Specific to our network, we utilize bottleneck to transform the features in vector form into a higher-dimensional space and project it to a hyperbolic manifold. This operation can be viewed as the immersion of the manifold space, which ensures that the high-dimensional manifold retains local properties of the manifold of low dimension, enabling the network to learn the geometric structure of the data.

IV. EXPERIMENTS AND ANALYSIS

This section describes the experiments in terms of data sets, the implementation, the training details, the comparison results, and the corresponding analysis.

A. Data Sets

We evaluate the performance of HMANet on four benchmark skeleton-based action recognition data sets. In all data sets, we use only the skeleton joint markers.

NTU RGB+D [36]: It is currently one of the largest 3-D action recognition data sets, containing RGB+D videos and skeleton data for human action recognition. The motion data was captured from 40 human objects by three Microsoft Kinect V2 cameras. There are 56 880 samples with four million frames in 60 categories, and the maximum number of frames in all samples is 300. Each body skeleton records 25 joints. The original benchmark provides two evaluation methods, namely, cross-subject (CS) and cross-view (CV) evaluation. In CS evaluation, the training set contains 40 320 videos from 20 subjects, and the remaining 16 560 videos are used for testing. In CV evaluation, 37 920 videos captured from No. 2 and No. 3 cameras were used for training, and the remaining 18 960 videos from No. 1 camera were used for testing. We follow the original two benchmarks and report the accuracy of Top-1.

Gaming-3D (G3D) [37]: It is a gaming data set collected with Microsoft Kinect which contains a total of 663 motion sequences. The data set consists of 20 actions performed by ten subjects in a controlled indoor environment. Each people performs several times and each sequence may contain multiple actions. As this data set consists of gaming actions, it has many temporal dependencies and rapid movements of body parts in the video sequences. The data set provides RGB video data and skeleton data. Skeleton data provides the 3-D coordinates of the joints. Each body in a sequence records 20 joints. We use the same protocol as the other works wherein the first five subjects are used for training, and the remaining for testing.

SHREC'17 Track Data Set [38]: The data set is a public dynamic hand gesture data set presented for the SHREC'17 Track. It contains sequences of 14 gestures performed between one and ten times by 28 participants in two-finger configurations, resulting in 2800 sequences. The data is categorized with two levels of granularity, presenting 14 and 28 actions, respectively. The coordinates of 22 hand joints in the 3-D world space are provided per frame, forming a full hand skeleton. Following the evaluation protocol of SHREC'17 track [38], we trained our model on 1960 samples and evaluated on the other 840 samples.

DHG-14/28 Data Set [38]: The data set is a public dynamic hand gesture data set collected by the Intel RealSense short-range depth camera. It contains sequences of 14 hand gestures performed time times by 20 participants, resulting in 2800 video sequences. The gestures are performed in two ways: 1) using one finger and 2) using the whole hand. The coordinates of 22 hand joints in the 3-D world space are provided per frame, forming a full hand skeleton. Although the DHG-14/28 data set has the same hand gestures with the SHREC'17 track data set, it is more challenging due to the leave-one-subject-out experimental protocol.

B. Implementation

Before the data was fed into the networks, we conduct some preprocessing such that the data structure of each video clip is unified. Since different actions last for various durations, the input sequences are normalized to a fixed length (128 for

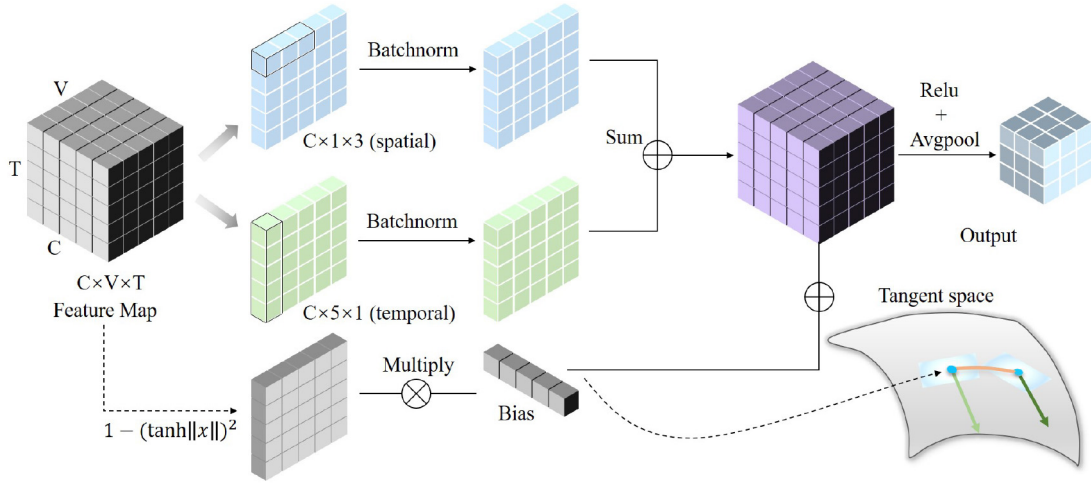


Fig. 3. Illustration of spatio-temporal convolution block equipped with manifold aware bias. Here, the convolution kernel is divided into the spatial kernel and temporal kernel and used on the same feature map, the outputs of the two are summed. Both of them are followed by a batch normalization (BN) layer. Moreover, a changeable bias based on a hyperbolic manifold is utilized to parallel transport the feature. The output is followed by activation layer (ReLU) and Avgpool is fed into the next block.

NTU RGB+D and 64 for others) through bilinear interpolation along the frame dimension. For the single-person sample in data set with two objects, the second body will be padded with all zeros. Each dimension of the 3-D coordinates is put into three channels as inputs. In order to evaluate the effectiveness of our framework more purely, we use relatively primitive data without any preprocessing such as random noise and random cropping.

The basic framework is as illustrated in Fig. 1. Take the NTU RGB+D data set as an example. First, we concatenate each dimension of the 3-D position tensor with the deviation tensor and designate joints as channels to extract the global response separately. Then, the position vectors and deviation vectors are mapped to the manifold along the coordinate dimension. They are summed on the manifold space to obtain the new position tensor. To conduct convolution in manifold space, we use the logarithmic map to project manifold features to Euclidean space. After the original data being transformed to a 128×128 pseudo image with six channels, a bottleneck is utilized to extend the feature to 64D. Then, the corresponding spatio-temporal Conv layers are built in tangent space. As shown in Fig. 3, in each convolution layer, 5×1 temporal convolution and 1×3 spatial convolution are coincident, and a bias adapted to manifold is used to move the features. We empirically stack six layers on this tangent space and channels at each layer are [64, 64, 128, 256, 256, 256]. Following each layer, a 2×2 Avgpooling is utilized to reduce dimensionality. Finally, the resulted features are averaged along the temporal dimension to four vectors in tangent space, and an FC layer followed by a softmax function is utilized to predict a class prediction.

During the training process, the cross-entropy loss is utilized as the classification loss. The learning rate is set as 0.01 and is decreased based on a cosine function. A stochastic gradient descent (SGD) with the Nesterov momentum (0.9) is applied as the optimization algorithm for the network. We set the weight decay to 0.0002 as regularization. This

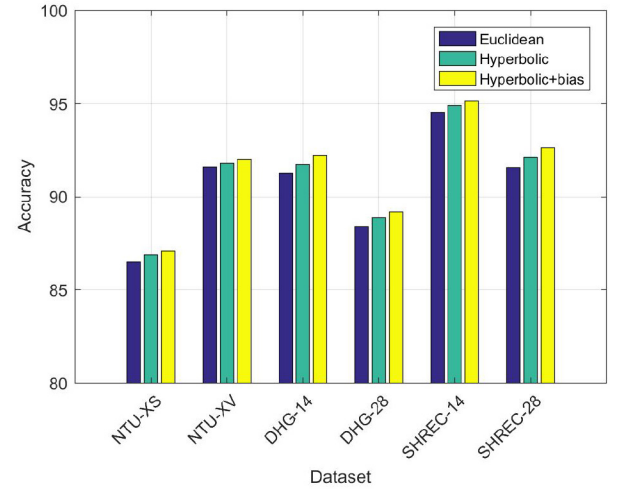


Fig. 4. Comparison of recognition accuracy of using and not using hyperbolic manifold aware mechanism on various data sets.

model will be trained for 70 epochs and compared to other approaches.

C. Ablation Study

We evaluate the effectiveness of our framework on the data sets of the human body and gesture data sets under given two evaluation measures. First, we evaluate how much benefit we obtained from hyperbolic geometry, we implemented our network without hyperbolic geometry. It works directly in the Euclidean space without Möbius addition with its corresponding distance tensor. Simultaneously, the changeable bias based on the hyperbolic manifold is also removed. This network serves as a baseline for comparison with our improved model. The comparison results are shown in Fig. 4. It can be seen from our experiments, with the help of hyperbolic space, for a given evaluation, our HMANet can improve the performance without increasing the parameters. Specifically, under the X-subject

TABLE II
AVERAGE VALUE OF THE STATISTICS OF SPHERE CENTER DISTANCE
UNDER TEN EXPERIMENTS

Statistics	Without Möbius addition	With Möbius addition
Minimum	0.5457	0.2169
Maximum	0.9486	0.9775
Variance	0.0075	0.0181

TABLE III
PERFORMANCE COMPARISON ON G3D USING CS PROTOCOL

Method	Year	Accuracy(%)
LRBM[39]	2015	90.50
R3D[27]	2016	90.94
JTM[9]	2018	94.24
CNN[9]	2018	96.00
HDM-BG[40]	2019	92.00
FIB-CNN[12]	2020	93.11
KM+TSC[41]	2021	92.91
Proposed HMANet	-	97.16

and X-view evaluations in NTU RGB+D, our model can overcome the baseline by 0.5% and 0.3%, respectively. Similarly, in SHREC'17, the proposed model defined in hyperbolic space could even outperform it by 1.1% with 28 gesture setting. All of them prove that defining the model on manifold space could benefit greatly.

To further evaluate the effectiveness of the proposed Hyperbolic aware bias, we remove the bias based on the proposed network and conduct comparative experiments. The result is shown by the green bar in the histogram. From the figure, we can notice that using the proposed bias has a certain accuracy improvement on all data sets, especially in gesture data. This is due to the fact that the gesture skeleton is extended from the wrist, which is more tree-like and has a clearer hierarchy between joints. To evaluate the influences of the Möbius addition on the distribution of features in hyperbolic space, we calculate the distance from a point in hyperbolic space to the center of the sphere in both cases, and use statistics to measure the sparsity of features, as shown in Table II. It can be seen from the table that after using the Möbius addition, the points are more dispersed in the space, and a hierarchical structure is more reflected, indicating that the network perceives the properties of the underlying hyperbolic manifold.

D. Comparison With State-of-the-Art

1) *G3D*: Table III shows comparison with previous methods. Our HMANet is able to achieve superior performance to [39] which extends the Restricted Boltzmann Machine. We are also able to outperform [27] which uses a rolling map to project data represented as points on Lie Group to the Euclidean space, and the recent method [40]. Work [9] encodes the 3-D skeleton data into 2-D images, and then utilizes the convolutional network for recognition, increasing the accuracy to 94.24%, which verifies the effectiveness of the

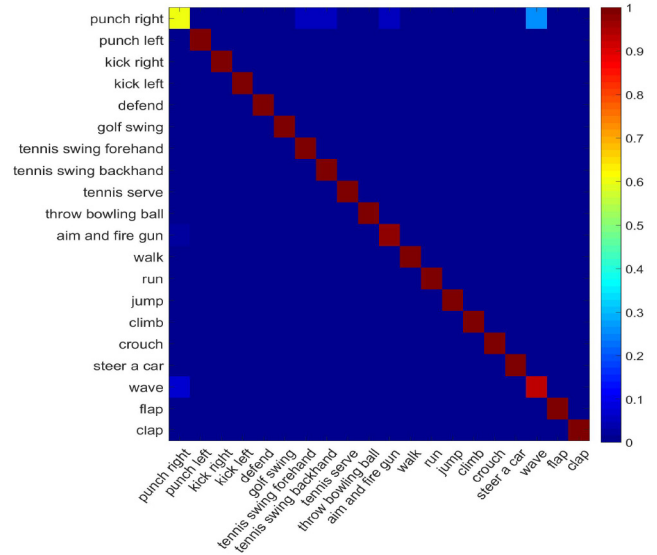


Fig. 5. Confusion matrix of G3D data set using CS protocol.

TABLE IV
PERFORMANCE COMPARISON ON NTU RGB+D USING
CS AND CV PROTOCOL

Method	Year	Cross-Subject(%)	Cross-View(%)
STA-LSTM[42]	2017	73.4	81.2
GCA-LSTM[43]	2017	74.4	82.8
DS-LSTM[44]	2020	77.80	87.33
CNN+LSTM[21]	2017	82.89	90.10
MTCNN+Rot.Clips[45]	2018	81.09	87.37
HCN[46]	2018	86.5	91.1
TSSI+GLAN+SSAN[20]	2019	82.4	89.1
(P+C)Net[19]	2019	86.1	93.5
PoF2I[22]	2019	82.46	89.53
TSRJI[47]	2019	73.3	80.3
POT2I+Inception v3[23]	2020	83.85	90.33
FIB-CNN[12]	2021	84.22	89.71
LAGA-Net[48]	2021	87.07	93.17
ST-GCN[13]	2018	81.5	88.3
GECNNs[49]	2020	85.4	91.1
Proposed HMANet	-	87.1	92.0

convolutional network. We are able to outperform other CNN-based methods [9], [12] largely, achieving a state-of-the-art result. It can be seen from the confusion matrix in Fig. 5, the recognition errors concentrate on punch right and wave, while our HMANet achieves almost 100% in all other recognizing actions.

2) *NTU RGB+D*: Table IV shows a comparison of our HMANet with past networks. When compared to the LSTM-based approaches [42]–[44], our network achieves superior performance. This is due to incorporating spatio-temporal features in our model which is sensitive to the geometrical information in the sequence. Furthermore, among all the CNN-based methods, we achieve the best result on CS protocol. Our HMANet outperforms approach [21] without employing LSTM networks and [20] without using a specific depth-first traversal. The CNN method [46] achieves the best performance after learning the tangent space features, which is also considered in our model. We have further achieved competitive performance to some GCN-based method [13], [49], and outperforming the pioneering work [13] on both CS and CV

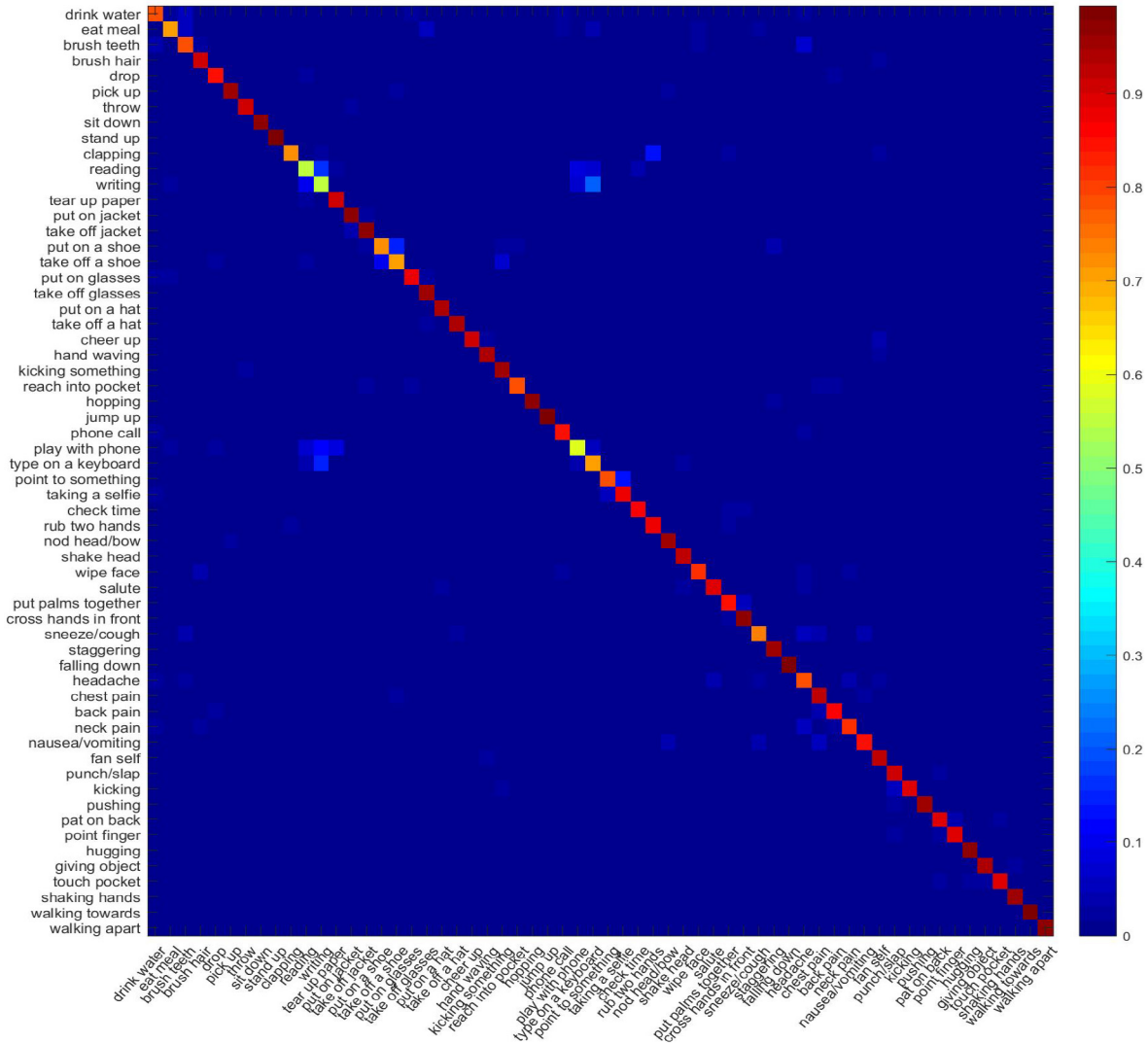


Fig. 6. Confusion matrix of NTU RGB+D data set in terms of the CS protocol.

protocol by 5.6% and 3.7%, respectively. When compared to the GCN-based methods, the CNN-based methods are unable to take full advantage of the topological structure due to a lack of input graph, hence do not perform well. The GCN-based methods utilize local information about specific body parts, while our network does not require any such separate handling of body parts. The confusion matrix of the classification results is shown in Fig. 6.

3) *SHREC'17 Track and DHG-14/28*: Table V shows the recognition accuracy of our framework trained and evaluated on SHREC'17 data set and DHG-14/28 data set. It shows that our HMANet achieves state-of-the-art performance under both 14 gesture and 28 gesture settings on the SHREC'17 data set, greater accuracy improvement with the more complicated 28 gestures setting, which further validates the effectiveness of our proposed model. Comparing the performance on DHG-14/28 data set, our proposed hierarchical architecture brings 6.61% and 8.08% accuracy improvement, respectively, for the 14 gestures setting and 28 gestures setting compared to CNN + LSTM. As shown in Table V, our network obtains 92.21% on 14 gesture protocol and 89.18%

on 28 gesture protocol. It is competitive with the state-of-the-art result in net HPEV+HMM [53] obtaining 92.54% and 88.86% for experiments with 14 and 28 gestures, respectively. Particularly, the good performance is more notable with 28 gestures setting than that with 14 gestures setting. Figs. 7 and 8 show the confusion matrix of our network on the DHG data set and the SHREC'17 data set. The recognition errors concentrate on highly similar actions, e.g., Grap to Pinch. Our network achieves 100% accuracy on both data sets in recognizing actions Swipe Right, Swipe+, and Swipe-V.

V. CONCLUSION

In this article, we propose a skeleton-based action recognition model using a hyperbolic manifold theory. The model is characterized by obtaining joint interaction from the spatial domain for the skeleton sequence, and parametrically representing it in hyperbolic space. Capturing the coordinate information of global joints through affine transformation, the spatial joint interaction can be fully explored to extract discriminative spatio-temporal features. Since the excellent ability

TABLE V
PERFORMANCE COMPARISON ON HAND GESTURE DATA SETS. 14G AND 28G REPRESENT 14 AND 28 GESTURE SETTINGS. (a) DHG-14/28 DATA SET USING THE LEAVE-ONE-SUBJECT-OUT PROTOCOL. (b) SHREC'17 TRACK DATA SET USING CS PROTOCOL

Method	Year	Accuracy(%)	
		14G	28G
CNN+LSTM[21]	2018	85.60	81.10
Res-TCN[50]	2018	86.90	83.60
STA-Res-TCN[50]	2018	89.20	85.00
ST-GCN[13]	2018	91.20	87.10
ST-TS-HGR-NET[51]	2019	87.30	83.40
SPD-NET[51]	2019	92.38	86.31
DG-STA[52]	2019	91.90	88.00
HPEV+HMM[53]	2020	92.54	88.86
Proposed HMANet	-	92.21	89.18

Method	Year	Accuracy(%)	
		14G	28G
CNN+LSTM[21]	2018	89.8	86.3
Parallel CNN[54]	2018	91.3	84.4
Res-TCN[50]	2018	91.1	87.3
STA-Res-TCN[50]	2018	93.6	90.7
MFA-Net[55]	2019	91.3	86.6
DD-Net[56]	2019	94.6	91.9
HPEV+HMM[53]	2020	94.88	92.26
TCN-Summ[57]	2021	93.57	91.43
Proposed HMANet	-	95.12	92.62

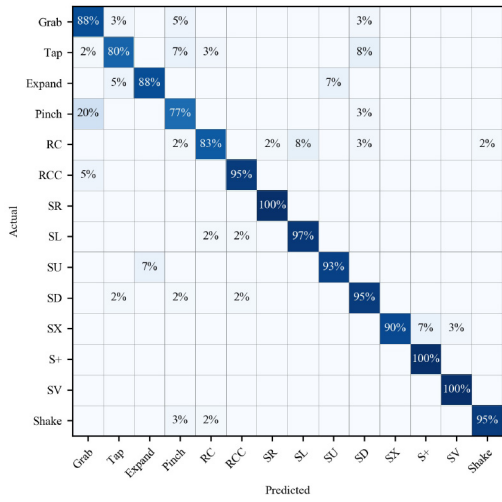


Fig. 7. Confusion matrix of DHG data set with 14 gestures setting.

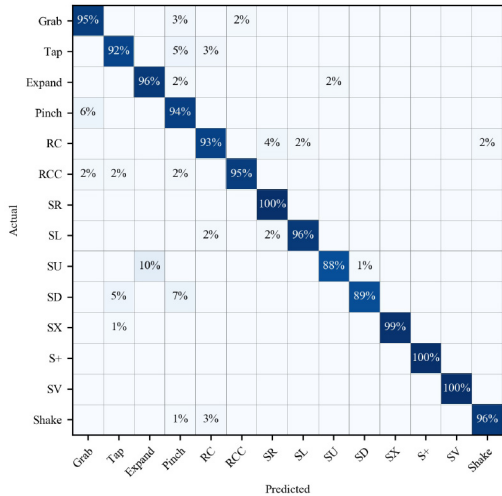


Fig. 8. Confusion matrix of SHREC'17 data set with 14 gestures setting.

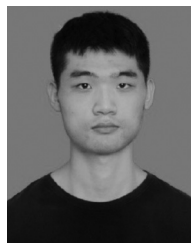
of hyperbolic space to represent hierarchical features, the global interactive features embedded can reflect certain hierarchical relationships. In addition, to extract spatio-temporal features, our method separately uses a temporal filter and

spatial filter to fuse local information and combine them, while a bias based on a hyperbolic manifold is added. On public data sets, including human body and gestures, we conduct experiments to prove that the proposed method has a certain versatility while achieving accuracy competitive with mainstream methods, which shows the method can be generalized to different skeleton structures. In the future, we will investigate how to fully embed skeleton features into hyperbolic space with lower distortion and lower computational complexity to better perceive the skeleton in hyperbolic space.

REFERENCES

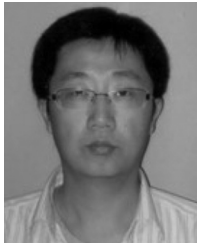
- [1] J. Shotton *et al.*, "Real-time human pose recognition in parts from single depth images," in *Proc. Int. Conf. Comput. Vis. Pattern Recognit.*, 2011, pp. 1297–1304.
- [2] D. Mehta *et al.*, "VNect: Real-time 3D human pose estimation with a single RGB camera," *ACM Trans. Graph.*, vol. 36, no. 4, pp. 1–14, 2017.
- [3] B. Sun, S. Wang, D. Kong, L. Wang, and B. Yin, "Real-time human action recognition using locally aggregated kinematic-guided Skeletonlet and supervised hashing-by-analysis model," *IEEE Trans. Cybern.*, early access, Sep. 7, 2021, doi: [10.48550/arXiv.2105.11312](https://doi.org/10.48550/arXiv.2105.11312).
- [4] Z. Shao, Y. Li, Y. Guo, X. Zhou, and S. Chen, "A hierarchical model for human action recognition from body-parts," *IEEE Trans. Circuits Syst. Video Technol.*, vol. 29, no. 10, pp. 2986–3000, Oct. 2019.
- [5] J. Weng, C. Weng, and J. Yuan, "Spatio-temporal naive-bayes nearest-neighbor (ST-NBNN) for skeleton-based action recognition," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit.*, 2017, pp. 445–454.
- [6] M. M. Arzani, M. Fathy, A. A. Azirani, and E. Adeli, "Switching structured prediction for simple and complex human activity recognition," *IEEE Trans. Cybern.*, vol. 51, no. 12, pp. 5859–5870, Dec. 2021.
- [7] S. Zhang, X. Liu, and J. Xiao, "On geometric features for skeleton-based action recognition using multilayer LSTM networks," in *Proc. Winter Conf. Appl. Comput. Vis. (WACV)*, 2017, pp. 148–157.
- [8] Q. Ma, Z. Chen, S. Tian, and W. W. Ng, "Difference-guided representation learning network for multivariate time-series classification," *IEEE Trans. Cybern.*, early access, Dec. 3, 2020, doi: [10.1109/TCYB.2020.3034755](https://doi.org/10.1109/TCYB.2020.3034755).
- [9] P. Wang, W. Li, C. Li, and Y. Hou, "Action recognition based on joint trajectory maps with convolutional neural networks," *Knowl.-Based Syst.*, vol. 158, pp. 43–53, Oct. 2018.
- [10] Y. Xu, J. Cheng, L. Wang, H. Xia, F. Liu, and D. Tao, "Ensemble one-dimensional convolution neural networks for skeleton-based action recognition," *IEEE Signal Process. Lett.*, vol. 25, no. 7, pp. 1044–1048, Jul. 2018.
- [11] P. Zhang, C. Lan, J. Xing, W. Zeng, J. Xue, and N. Zheng, "View adaptive neural networks for high performance skeleton-based human action recognition," *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 41, no. 8, pp. 1963–1978, Aug. 2019.

- [12] A. Banerjee, P. K. Singh, and R. Sarkar, "Fuzzy integral based CNN classifier fusion for 3D skeleton action recognition," *IEEE Trans. Circuits Syst. Video Technol.*, vol. 31, no. 6, pp. 2206–2216, Jun. 2021.
- [13] S. Yan, Y. Xiong, and D. Lin, "Spatial temporal graph convolutional networks for skeleton-based action recognition," in *Proc. AAAI Conf. Artif. Intell.*, 2018, pp. 7444–7452.
- [14] B. Li, X. Li, Z. Zhang, and F. Wu, "Spatio-temporal graph routing for skeleton-based action recognition," in *Proc. AAAI Conf. Artif. Intell.*, vol. 33, 2019, pp. 8561–8568.
- [15] F. Sala, C. De Sa, A. Gu, and C. Ré, "Representation tradeoffs for hyperbolic embeddings," in *Proc. Int. Conf. Mach. Learn. (PMLR)*, 2018, pp. 4457–4466.
- [16] R. Benedetti and C. Petronio, *Lectures on Hyperbolic Geometry*. Cham, Switzerland: Springer, 1992.
- [17] C. Gulcehre *et al.*, "Hyperbolic attention networks," 2018, *arXiv:1805.09786*.
- [18] C. Li, Y. Hou, P. Wang, and W. Li, "Joint distance maps based action recognition with convolutional neural networks," *IEEE Signal Process. Lett.*, vol. 24, no. 5, pp. 624–628, May 2017.
- [19] C. Cao, C. Lan, Y. Zhang, W. Zeng, H. Lu, and Y. Zhang, "Skeleton-based action recognition with gated convolutional neural networks," *IEEE Trans. Circuits Syst. Video Technol.*, vol. 29, no. 11, pp. 3247–3257, Nov. 2019.
- [20] Z. Yang, Y. Li, J. Yang, and J. Luo, "Action recognition with Spatio-Temporal visual attention on skeleton image sequences," *IEEE Trans. Circuits Syst. Video Technol.*, vol. 29, no. 8, pp. 2405–2415, Aug. 2019.
- [21] C. Li, P. Wang, S. Wang, Y. Hou, and W. Li, "Skeleton-based action recognition using LSTM and CNN," in *Proc. Int. Conf. Multimedia Expo Workshops (ICMEW)*, 2017, pp. 585–590.
- [22] T. Huynh-The, C.-H. Hua, and D.-S. Kim, "Encoding pose features to images with data augmentation for 3-D action recognition," *IEEE Trans. Ind. Informat.*, vol. 16, no. 5, pp. 3100–3111, May 2020.
- [23] T. Huynh-The, C.-H. Hua, T.-T. Ngo, and D.-S. Kim, "Image representation of pose-transition feature for 3D skeleton-based action recognition," *Inf. Sci.*, vol. 513, pp. 112–126, Mar. 2020.
- [24] T. Huynh-The, C.-H. Hua, N. A. Tu, and D.-S. Kim, "Learning geometric features with dual-stream CNN for 3D action recognition," in *Proc. Int. Conf. Acoust. Speech Signal Process. (ICASSP)*, 2020, pp. 2353–2357.
- [25] R. Vemulapalli, F. Arrate, and R. Chellappa, "Human action recognition by representing 3D skeletons as points in a lie group," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit.*, 2014, pp. 588–595.
- [26] Z. Huang, C. Wan, T. Probst, and L. Van Gool, "Deep learning on lie groups for Skeleton-based action recognition," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit.*, 2017, pp. 6099–6108.
- [27] R. Vemulapalli and R. Chellappa, "Rolling rotations for recognizing human actions from 3D Skeletal data," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit.*, 2016, pp. 4471–4479.
- [28] E. Mathieu, C. L. Lan, C. J. Maddison, R. Tomioka, and Y. W. Teh, "Continuous hierarchical representations with Poincaré variational auto-encoders," 2019, *arXiv:1901.06033*.
- [29] O.-E. Ganea, G. Bécigneul, and T. Hofmann, "Hyperbolic neural networks," 2018, *arXiv:1805.09112*.
- [30] A. Ungar, "A Gyrovector space approach to hyperbolic geometry," *Synthesis Lectures Math. Statist.*, vol. 1, no. 1, pp. 1–194, 2008.
- [31] M. Nickel and D. Kiela, "Poincaré embeddings for learning hierarchical representations," in *Proc. Int. Conf. Neural Inf. Process. Syst.*, vol. 30, 2017, pp. 6338–6347.
- [32] A. Tifrea, G. Bécigneul, and O.-E. Ganea, "Poincaré glove: Hyperbolic word embeddings," 2018, *arXiv preprint arXiv:1810.06546*.
- [33] I. Chami, Z. Ying, C. Ré, and J. Leskovec, "Hyperbolic graph convolutional neural networks," in *Proc. Int. Conf. Neural Inf. Process. Syst.*, vol. 32, 2019, pp. 4868–4879.
- [34] M. B. Hauser, "Principles of Riemannian geometry in neural networks," in *Proc. NIPS*, 2018, pp. 2804–2813.
- [35] J. W. Cannon, W. J. Floyd, R. Kenyon, and W. R. Parry, "Hyperbolic geometry," *Flavors Geom.*, vol. 31, nos. 59–115, p. 2, 1997.
- [36] A. Shahroudy, J. Liu, T.-T. Ng, and G. Wang, "NTU RGB+ D: A large scale dataset for 3d human activity analysis," in *Proc. Int. Conf. Comput. Vis. Pattern Recognit.*, 2016, pp. 1010–1019.
- [37] V. Bloom, D. Makris, and V. Argyriou, "G3D: A gaming action dataset and real time action recognition evaluation framework," in *Proc. Int. Conf. Comput. Vis. Pattern Recognit.*, 2012, pp. 7–12.
- [38] Q. De Smedt, H. Wannous, J.-P. Vandeboer, J. Guerry, B. Le Saux, and D. Filliat, "Shrec'17 track: 3D hand gesture recognition using a depth and skeletal dataset," in *Proc. 3DOR-10th Eurograph. Workshop 3D Object Retrieval*, 2017, pp. 1–6.
- [39] S. Nie, Z. Wang, and Q. Ji, "A generative restricted Boltzmann machine based method for high-dimensional motion data modeling," *Comput. Vis. Image Understand.*, vol. 136, pp. 14–22, Jul. 2015.
- [40] R. Zhao, W. Xu, H. Su, and Q. Ji, "Bayesian hierarchical dynamic model for human action recognition," in *Proc. Int. Conf. Comput. Vis. Pattern Recognit.*, 2019, pp. 7733–7742.
- [41] G. Paoletti, J. Cavazza, C. Beyan, and A. D. Bue, "Subspace clustering for action recognition with covariance representations and temporal pruning," in *Proc. 25th Int. Conf. Pattern Recognit. (ICPR)*, 2021, pp. 6035–6042.
- [42] S. Song, C. Lan, J. Xing, W. Zeng, and J. Liu, "An end-to-end spatio-temporal attention model for human action recognition from skeleton data," in *Proc. AAAI Conf. Artif. Intell.*, vol. 31, 2017, pp. 4263–4270.
- [43] J. Liu, G. Wang, P. Hu, L.-Y. Duan, and A. C. Kot, "Global context-aware attention lstm networks for 3D action recognition," in *Proc. Int. Conf. Comput. Vis. Pattern Recognit.*, 2017, pp. 1647–1656.
- [44] X. Jiang, K. Xu, and T. Sun, "Action recognition scheme based on skeleton representation with DS-LSTM network," *IEEE Trans. Circuits Syst. Video Technol.*, vol. 30, no. 7, pp. 2129–2140, Jul. 2020.
- [45] Q. Ke, M. Bennamoun, S. An, F. Soheli, and F. Boussaid, "Learning clip representations for skeleton-based 3D action recognition," *IEEE Trans. Image Process.*, vol. 27, no. 6, pp. 2842–2855, Jun. 2018.
- [46] C. Li, Q. Zhong, D. Xie, and S. Pu, "Co-occurrence feature learning from skeleton data for action recognition and detection with hierarchical aggregation," 2018, *arXiv:1804.06055*.
- [47] C. Caetano, F. Brémont, and W. R. Schwartz, "Skeleton image representation for 3d action recognition based on tree structure and reference joints," in *Proc. 32nd SIBGRAPI Conf. Graph. Patterns Images (SIBGRAPI)*, 2019, pp. 16–23.
- [48] R. Xia, Y. Li, and W. Luo, "LAGA-net: Local-and-global attention network for skeleton based action recognition," *IEEE Trans. Multimedia*, early access, Jun. 7, 2021, doi: [10.1109/TMM.2021.3086758](https://doi.org/10.1109/TMM.2021.3086758).
- [49] X. Zhang, C. Xu, X. Tian, and D. Tao, "Graph edge convolutional neural networks for Skeleton-based action recognition," *IEEE Trans. Neural Netw. Learn. Syst.*, vol. 31, no. 8, pp. 3047–3060, Aug. 2020.
- [50] J. Hou, G. Wang, X. Chen, J.-H. Xue, R. Zhu, and H. Yang, "Spatial-temporal attention res-TCN for Skeleton-based dynamic hand gesture recognition," in *Proc. Eur. Conf. Comput. Vis. (ECCV) Workshops*, 2018, pp. 273–286.
- [51] X. S. Nguyen, L. Brun, O. Lézoray, and S. Bougleux, "A neural network based on SPD manifold learning for Skeleton-based hand gesture recognition," in *Proc. Int. Conf. Comput. Vis. Pattern Recognit.*, 2019, pp. 12036–12045.
- [52] Y. Chen, L. Zhao, X. Peng, J. Yuan, and D. N. Metaxas, "Construct dynamic graphs for hand gesture recognition via spatial-temporal attention," 2019, *arXiv:1907.08871*.
- [53] J. Liu, Y. Liu, Y. Wang, V. Prinnet, S. Xiang, and C. Pan, "Decoupled representation learning for Skeleton-based gesture recognition," in *Proc. Int. Conf. Comput. Vis. Pattern Recognit.*, 2020, pp. 5750–5759.
- [54] G. Devineau, W. Xi, F. Moutarde, and J. Yang, "Convolutional neural networks for multivariate time series classification using both inter- and intra-channel parallel convolutions," in *Proc. Reconnaissance des Formes Image Apprentissage et Perception (RFIAP)*, 2018.
- [55] X. Chen, H. Guo, G. Wang, and L. Zhang, "Motion feature augmented recurrent neural network for skeleton-based dynamic hand gesture recognition," in *Proc. Int. Conf. Image Process. (ICIP)*, 2017, pp. 2881–2885.
- [56] F. Yang, Y. Wu, S. Sakti, and S. Nakamura, "Make skeleton-based action recognition model smaller, faster and better," in *Proc. ACM Multimedia Asia*, 2019, pp. 1–6.
- [57] A. Sabater, I. Alonso, L. Montesano, and A. C. Murillo, "Domain and view-point agnostic hand action recognition," 2021, *arXiv:2103.02303*.



Jinghong Chen is currently pursuing the post-graduate degree with the Shenzhen Research Institute, Xiamen University, Shenzhen, China, and also with the Department of Computer Science and Technology, Xiamen University, Xiamen, China.

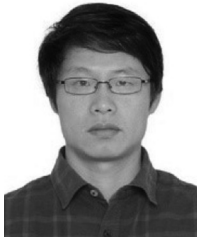
His research interests include action recognition, computer vision, and deep learning.



Chong Zhao received the B.S. degree from the Department of Computer Science, Jilin University, Changchun, China, the M.S. degree from the Academy of Mathematics and Systems Science, Chinese Academy of Sciences, Beijing, China, and the Ph.D. degree from the Department of Computer Science and Engineering, The Chinese University of Hong Kong, Hong Kong.

He is currently an Assistant Professor with the Department of Computer Science, Xiamen University, Xiamen, China. His research interests

include geometry processing, computer graphics, and computer vision.



Qicong Wang received the Ph.D. degree from Zhejiang University, Hangzhou, China.

He is currently an Associate Professor with the Shenzhen Research Institute, Xiamen University, Shenzhen, China, and also with the Department of Computer Science and Technology, Xiamen University, Xiamen, China. His research interests include machine vision, robot navigation, and machine learning.



Hongying Meng (Senior Member, IEEE) received the Ph.D. degree in communication and electronic systems from Xi'an Jiaotong University, Xi'an China.

He has authored over 130 publications, including IEEE TRANSACTIONS ON IMAGE PROCESSING, IEEE TRANSACTIONS ON CYBERNETICS, IEEE TRANSACTIONS ON FUZZY SYSTEMS, IEEE TRANSACTIONS ON AUTOMATIC CONTROL, IEEE TRANSACTIONS ON CIRCUITS AND SYSTEMS FOR VIDEO TECHNOLOGY

(TCSVT), IEEE TRANSACTIONS ON BIOMEDICAL ENGINEERING, IEEE TRANSACTIONS ON COGNITIVE AND DEVELOPMENTAL SYSTEMS (TCDS), ICASSP, and CVPR. His research interests include digital signal processing, affective computing, machine learning, human-computer interaction, and computer vision.

Dr. Meng is an Associate Editor for TCSVT and TCDS. He is currently a Reader with the Department of Electronic and Computer Engineering, Brunel University London, London, U.K. He is a Fellow of The Higher Education Academy and a member of Engineering Professors Council in U.K.