

Real-time 3D human action recognition based on Hyperpoint sequence

Xing Li, *Student Member, IEEE*, Qian Huang, Zhijian Wang, Tianjin Yang, Zhenjie Hou, and Zhuang Miao

Abstract—Real-time 3D human action recognition has broad industrial applications, such as surveillance, human-computer interaction, and healthcare monitoring. By relying on complex spatio-temporal local encoding, most existing point cloud sequence networks capture spatio-temporal local structures to recognize 3D human actions. To simplify the point cloud sequence modeling task, we propose a lightweight and effective point cloud sequence network referred to as SequentialPointNet for real-time 3D action recognition. Instead of capturing spatio-temporal local structures, SequentialPointNet encodes the temporal evolution of static appearances to recognize human actions. Firstly, we define a novel type of point data, Hyperpoint, to better describe the temporally changing human appearances. A theoretical foundation is provided to clarify the information equivalence property for converting point cloud sequences into Hyperpoint sequences. Secondly, the point cloud sequence modeling task is decomposed into a Hyperpoint embedding task and a Hyperpoint sequence modeling task. Specifically, for Hyperpoint embedding, the static point cloud technology is employed to convert point cloud sequences into Hyperpoint sequences, which introduces inherent frame-level parallelism; for Hyperpoint sequence modeling, a Hyperpoint-Mixer module is designed as the basic building block to learning the spatio-temporal features of human actions. Extensive experiments on three widely-used 3D action recognition datasets demonstrate that the proposed SequentialPointNet achieves competitive classification performance with up to 10X faster than existing approaches.

Index Terms—3D action recognition, point cloud sequence, SequentialPointNet, Hyperpoint.

I. INTRODUCTION

WITH the release of low-cost depth cameras, 3D human action recognition has attracted more and more attention from researchers, which has broad industrial applications, such as human-computer interaction, security surveillance, automated driving, and robotics applications [1]. 3D human action recognition [2]–[4] can be divided into three categories based on different data types: 1) skeleton sequence-

Manuscript received May 6, 2022. This work is partly supported by the National Key Research and Development Program of China under grant No. 2018YFC0407905. (Corresponding authors: Qian Huang.)

Xing Li, Qian Huang, Zhijian Wang, and Tianjin Yang are with the Key Laboratory of Water Big Data Technology of Ministry of Water Resources, Hohai University; and School of Computer and Information, Hohai University, Nanjing 211100, China (e-mail: huangqian@hhu.edu.cn).

Zhenjie Hou is with Changzhou University, Changzhou 213000, China.

Zhuang Miao is with Command and Control Engineering College, Army Engineering University of PLA, Nanjing 210007, China (e-mail: emiao.beyond@163.com).

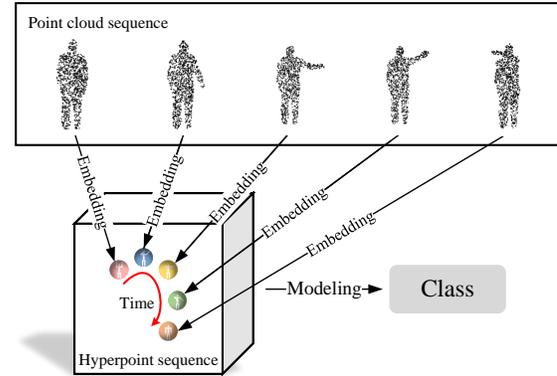


Fig. 1. SequentialPointNet decomposes the complex point cloud sequence modeling task into a static point cloud technology-based Hyperpoint embedding task and a simple Hyperpoint sequence modeling task.

based 3D human action recognition, which focuses on human joint trajectories; 2) depth sequence-based 3D human action recognition, in which pixel values in depth maps represent the distance information from human bodies to the depth camera; 3) point cloud sequence-based 3D human action recognition, which studies the spatial distribution of human appearance along the time dimension. Compared with skeleton sequence-based approaches, point cloud sequences are more convenient to be collected without additional pose estimation algorithms. Compared with depth sequence-based approaches, point cloud sequence-based methods yield lower computation costs. For these reasons, we focus on point cloud sequence-based 3D human action recognition in this work.

Point cloud sequences of 3D human actions exhibit unordered intra-frame spatial information and ordered inter-frame temporal information. Due to the complex data structures, capturing the spatio-temporal textures from point cloud sequences is extremely challenging. An intuitive way is to convert the point cloud sequence to a static point cloud and employ the static point cloud method to encode [5]. However, using static point clouds to represent the entire point cloud sequence loses a lot of spatio-temporal information, which degrades the recognition performance. Therefore, point cloud sequence methods are necessary, which directly consume the point cloud sequence for 3D human action classification. In order to model the dynamics of point cloud sequences, cross-frame spatio-temporal local neighborhoods are usually constructed. After that, point spatio-temporal operations, such as 4D PointNet [6], point spatio-temporal convolution [7],

and point 4D convolution [8], are used to encode the spatio-temporal local structures. Nevertheless, unlike conventional grid-based videos, spatio-temporal local encoding for point cloud sequences is both sophisticated and time-consuming because point cloud sequences are irregular and unordered in the spatial dimension. Furthermore, cross-frame spatio-temporal local encoding restricts frame-level parallel computation, which is not conducive to real-time point cloud sequence modeling.

In this paper, we propose a lightweight frame-level parallel point cloud sequence network, named SequentialPointNet, for real-time 3D human action recognition without resorting to spatio-temporal local encoding. SequentialPointNet remarkably simplifies the complexity of the point cloud sequence modeling task and improves computational efficiency while achieving strong recognition performance. SequentialPointNet encodes the temporal evolution of static appearances instead of capturing spatio-temporal local structures to recognize human actions. We define a novel type of point data named Hyperpoint to better describe the temporally changing human appearances. As shown in Fig. 1, the key to our approach is to decompose the complex point cloud sequence modeling task into a static point cloud technology-based Hyperpoint embedding task and a simple Hyperpoint sequence modeling task. Specifically, we employ the static point cloud technology as a Hyperpoint embedding module to flatten each point cloud frame into a Hyperpoint. In this fashion, the spatial structure of the human appearance is encoded and embedded into the corresponding Hyperpoint. Further, these Hyperpoints are assembled together and form a Hyperpoint sequence. We design a Hyperpoint-Mixer module as the basic building block to model Hyperpoint sequences for recognizing 3D human actions.

Our main contributions are summarized as follows:

- To avoid computationally expensive spatio-temporal local encoding, we propose an efficient and effective point cloud sequence network, dubbed SequentialPointNet, for real-time 3D human action recognition. SequentialPointNet treats point cloud sequence modeling as a two-phase task: Hyperpoint embedding and Hyperpoint sequence modeling, which significantly simplifies the encoding complexity and improves computational efficiency while realizing superior recognition performance.

- To the best of our knowledge, we are the first to mathematically define Hyperpoint and design a Hyperpoint-Mixer module as the basic building block to model Hyperpoint sequences.

- Our SequentialPointNet achieves up to $10\times$ faster than existing point cloud sequence models and yields the cross-view accuracy of 97.6% on the NTU RGB+D 60 dataset, the cross-setup accuracy of 95.4% on the NTU RGB+D 120 dataset, and the cross-subject accuracy of 92.64% on the MSR Action3D dataset, which outperforms the state-of-the-art methods. Moreover, our SequentialPointNet adapted to skeleton sequences achieves better performance compared to the most famous skeleton sequence-based action recognition method.

This paper is an extension of our previous work Hyper-

pointNet [9]. In HyperpointNet, to avoid spatio-temporal local encoding, the Hyperpoint is proposed as a new point data type but its mathematical definition is not provided. HyperpointNet converts the point cloud sequence into a Hyperpoint sequence and further converts the Hyperpoint sequence into static point clouds. By doing so, HyperpointNet has the ability to model point cloud sequences utilizing only static point cloud encoding techniques. In SequentialPointNet, we provide a detailed mathematical definition and application scenarios of Hyperpoint sequence and compare it with existing point cloud data. SequentialPointNet explicitly decomposes the point cloud sequence modeling task into a Hyperpoint embedding task and a Hyperpoint sequence modeling task, as well as gives a theoretical foundation for transforming point cloud sequences into Hyperpoint sequences. For the Hyperpoint sequence modeling task, SequentialPointNet designs a specialized network architecture named Hyperpoint-Mixer for Hyperpoint sequences rather than transforming them into static point clouds. Experiments also demonstrate that SequentialPointNet is superior to HyperpointNet. Compared to our previous work, we include another dataset for 3D action recognition, *i.e.*, MSR Action3D dataset [10]. Moreover, we further apply our SequentialPointNet to skeleton sequence-based 3D human action recognition on NTU RGB+D 60 dataset [11], which verifies the effectiveness of SequentialPointNet's adaption to a skeleton sequence modeling task.

SequentialPointNet's code has been made public at <https://github.com/XingLi1012/SequentialPointNet.git>.

II. RELATED WORK

A. Static Point Cloud Modeling

With the popularity of low-cost 3D sensors, deep learning on static point clouds has attracted much attention from researchers due to extensive applications ranging from object classification, part segmentation, to scene semantic parsing. Static point clouds models [12], [13] can be divided into volumetric-based methods and point-based methods. Volumetric-based methods usually voxelize a point cloud into 3D grids, and then a 3D Convolution Neural Network (CNN) is applied on the volumetric representation for classification. Point-based methods are directly performed on raw point clouds. PointNet [13] is a pioneering effort that directly processes point sets. The key idea of PointNet is to abstract each point using a set of Multi-layer Perceptrons (MLPs) and then assemble all individual point features by a symmetry function, *i.e.*, a max pooling operation. However, PointNet lacks the ability to capture local structures. Therefore, in [14], a hierarchical network PointNet++ is proposed to encode fine geometric structures from the neighborhood of each point. PointNet++ is made of several set abstraction levels. A set abstraction level is composed of three layers: sampling layer, grouping layer, and PointNet-based learning layer. By stacking several set abstraction levels, PointNet++ progressively abstracts larger and larger local regions along the hierarchy.

B. Point Cloud Sequence-based 3D Human Action Recognition

Point cloud sequence-based 3D human action recognition is a fairly new and challenging task. 3DV-PointNet++ [5] is the first volumetric-based work to recognize human actions from point cloud sequences. In 3DV-PointNet++, 3D dynamic voxel (3DV) is proposed as a novel 3D motion representation. A set of points is extracted from 3DV and input into PointNet++ for 3D action recognition in the end-to-end learning way. However, since the point cloud sequence is converted into a static 3D point cloud set, 3DV-PointNet++ loses a lot of spatio-temporal information and increases additional computational costs.

To overcome this problem, researchers have focused mainly on investigating point cloud sequence networks that directly consume point cloud sequences for human action recognition. MeteorNet [6] is the first work on deep learning for modeling point cloud sequences. In MeteorNet, two ways are proposed to construct spatio-temporal local neighborhoods for each point in the point cloud sequence. The abstracted features of each point are learned by aggregating the information from these neighborhoods. Fan *et al.* [7] propose a Point spatio-temporal (PST) convolution to encode the spatio-temporal local structures of point cloud sequences. PST convolution first disentangles space and time in point cloud sequences. Then, PST convolutions are incorporated into a deep network namely PSTNet to model point cloud sequences in a hierarchical manner. To avoid point tracking, Point 4D Transformer (P4Transformer) network [8] is proposed to model point cloud videos. Specifically, P4Transformer consists of a point 4D convolution to embed the spatio-temporal local structures presented in a point cloud video and a Transformer to encode the appearance and motion information by performing self-attention on the embedded local features. However, in these point cloud sequence networks, spatio-temporal local encoding is usually performed during modeling point cloud sequences, which is quite time-consuming and limits the real-time computational efficiency. In HyperpointNet [9], the point cloud sequence is first abstracted into a new point data type named Hyperpoint sequence. Then, the Hyperpoint sequence is converted into static point clouds by injecting order information. In this fashion, the point cloud sequence modeling task is decomposed into two static point cloud encoding tasks. HyperpointNet is the first deep neural network to model point cloud sequences without performing cross-frame spatio-temporal local encoding. However, HyperpointNet, only based on static point cloud technology, does not fully exploit the spatio-temporal information of the Hyperpoint sequence.

III. METHODOLOGY

In this section, we present a lightweight and effective point cloud sequence network referred to as SequentialPointNet for real-time 3D action recognition. Rather than capturing spatio-temporal local structures, our SequentialPointNet encodes the temporal evolution of static appearances to recognize human actions. By innovatively decomposing the modeling task of point cloud sequences into a Hyperpoint embedding task

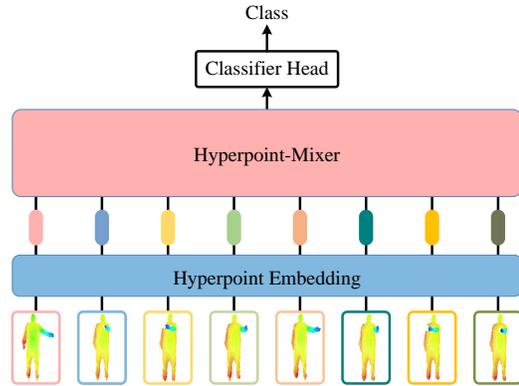


Fig. 2. SequentialPointNet contains a Hyperpoint embedding module, a Hyperpoint-Mixer module, and a classifier head.

and a Hyperpoint sequence modeling task, SequentialPointNet significantly reduces the complexity and time cost for point cloud sequence modeling, while attaining very high recognition accuracy. The overall flowchart of SequentialPointNet is described in Fig. 2. SequentialPointNet contains a Hyperpoint embedding module, a Hyperpoint-Mixer module, and a classifier head.

A. Hyperpoint Sequence

We define a new type of point data named Hyperpoint. A Hyperpoint sequence of length T is represented as $\mathbf{F} = \{\mathbf{F}_t \mid t = 1, \dots, T\}$ with $\mathbf{F}_t \in \mathbb{R}^{d_H+m_H}$, which is an ordered set of high-dimensional points. Each Hyperpoint \mathbf{F}_t is a vector of its d_H -dim coordinate plus extra m_H -dim feature channels. Extra feature channels can be neighborhood structure information or other modal features from different sensors. In this paper, unless otherwise noted, we only use the high-dimensional coordinate as our Hyperpoint's channels.

The Hyperpoint sequence can be used to describe a temporally changing complex information unit, which exists widely in the real world such as human spatio-temporal actions, 4D driving scenarios, etc. In this paper, the Hyperpoint sequence is obtained by flattening each frame of the point cloud sequence. Each Hyperpoint encapsulates the complex spatial structure of the human appearance at a specific moment.

Comparison with static point cloud and point cloud sequence: The static point cloud or individual point cloud frame is essentially an unordered point set. Ordered point cloud frames further form a point cloud sequence. Point cloud sequences are ordered in the temporal dimension but spatially irregular. In contrast, Hyperpoint sequences do not include unordered point sets and are easier to be encoded. Furthermore, unlike static point clouds and point cloud sequences, each element in the Hyperpoint sequence carries complex internal information. The comparison of the three types of point data is presented in Fig. 3.

B. Hyperpoint Embedding Module

To extract features of video data, spatio-temporal local encoding is usually performed. However, as point cloud se-

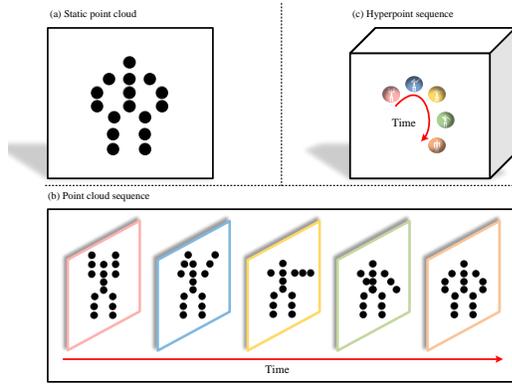


Fig. 3. The comparison of the three types of point data.

quences are irregular and lack order along the spatial dimension, spatio-temporal local encoding for point cloud sequences is complex and time-consuming. Moreover, spatio-temporal local encoding limits real-time point cloud sequence classification due to heavy cross-frame computational dependencies.

To avoid spatio-temporal local encoding, SequentialPointNet disentangles space and time in point cloud sequences. In our SequentialPointNet, the Hyperpoint embedding module is performed to preferentially encode each point cloud frame and output a Hyperpoint summarizing spatial structures. After that, the complex point cloud sequence modeling task is transformed into a simple Hyperpoint sequence modeling task. An outstanding advantage of this space-first strategy is that existing static point cloud models [13], [14] can be used almost out of the box as the Hyperpoint embedding module. In addition, embedding each point cloud frame independently allows the Hyperpoint sequences to be generated in a frame-level parallel manner, thereby improving modeling efficiency.

In the Hyperpoint embedding module, each point cloud frame is transformed into a Hyperpoint summarizing human static appearance by the static point cloud technology. We first adopt the set abstraction operation twice to downsample each point cloud frame. In this process, the texture information from space partitions is aggregated into the corresponding centroids. Then, in order to characterize human entire appearance information, a PointNet layer is used.

The set abstraction operation in our work is made of three key layers: sampling layer, grouping layer, and augmentation PointNet layer. Specifically, let $\mathcal{S} = \{\mathcal{S}_t\}_{t=1}^T$ denote a point cloud sequence of T frames, and $\mathcal{S}_t = \{x_1^t, x_2^t, \dots, x_n^t\}$ denotes the unordered point set of the t -th frame, where n is the number of points. The m -th set abstraction level takes an $n_{m-1} \times (d + c_{m-1})$ matrix as input that is from n_{m-1} points with d -dim coordinates and c_{m-1} -dim point features. It outputs an $n_m \times (d + c_m)$ matrix of n_m downsampled points with d -dim coordinates and new c_m -dim feature vectors summarizing local context. The size of the input in the first set abstraction level is $n \times (d + c)$. In this work, d is set as 3 corresponding to the 3D coordinate (X, Y, Z) of each point, and c is set as 0.

In the sampling layer, farthest point sampling (FPS) [14] is used to choose n_m points as centroids from the point set.

In the grouping layer, a point set of size $n_{m-1} \times (d + c_{m-1})$ and the coordinates of a set of centroids of size $n_m \times d$ are taken as input. The output is n_m groups of point sets of size $n_m \times k_m \times (d + c_{m-1})$, where each group corresponds to a local region and k_m is the number of local points in the neighborhood of centroid points. Ball query finds all points that are within a radius to the query point, in which an upper limit of k_m is set.

The augmentation PointNet layer in the set abstraction operation includes an inter-feature attention mechanism, a set of MLPs, and a max pooling operation. The input of this layer is n_m local regions with data size $n_m \times k_m \times (d + c_{m-1})$. First, the coordinates of points in a local region are translated into a local frame relative to the centroid point. Second, the distance between each local point and the corresponding centroid is used as a 1D additional point feature to alleviate the influence of rotational motion on action recognition. Then, an inter-feature attention mechanism is adopted to optimize the fusion effect of different features. The inter-feature attention mechanism is realized by Channel Attention Module (CAM). The inter-feature attention mechanism is not used in the first set abstraction operation due to only the 1-dim point feature. In the following, a set of MLPs are applied to abstract the features of each local point. Further, the representation of a local region is generated by incorporating the abstracted features of all local points using a max pooling operation. Finally, coordinates of the centroid point and its local region representation are concatenated as abstracted features of this centroid point. The augmentation PointNet layer is formalized as follows:

$$r_j^t = \left[\text{MAX}_{i=1, \dots, k_m} \left\{ \text{MLP} \left(\left[(l_{j,i}^t - o_j^t); e_{j,i}^t; p_{j,i}^t \right] \odot A \right) \right\}; o_j^t \right] \quad (1)$$

where $l_{j,i}^t$ is the coordinates of i -th point in the j -th local region from the t -th point cloud frame. o_j^t and $p_{j,i}^t$ are the coordinates of the centroid point and the point features corresponding to $l_{j,i}^t$, respectively. $e_{j,i}^t$ is the euclidean distance between $l_{j,i}^t$ and o_j^t . A denotes the attention mechanism with $(3+1+c_{m-1})$ -dim scores corresponding to the coordinates and features of each point. Attention scores in A are shared among all local points from all point cloud frames. \odot indicates dot product operation. r_j^t is the abstracted features of the j -th centroid point from the t -th point cloud frame.

The set abstract operation is performed twice in the Hyperpoint sequence module. In order to characterize the spatial appearance information of the entire point cloud frame, a PointNet layer consisting of a set of MLPs and a max pooling operation is used as follows:

$$F_t = \text{MAX}_{j=1, \dots, n_2} \left\{ \text{MLP} (r_j^t) \right\} \quad (2)$$

where F_t is the Hyperpoint of the t -th point cloud frame. So the Hyperpoint sequence is represented as $F = \{F_t\}_{t=1}^T$.

Theoretical foundation for converting point cloud sequences to Hyperpoint sequences: Temporally changing static appearances constitute the human spatio-temporal action. Our SequentialPointNet encodes the temporal evolution of static appearances instead of capturing spatio-temporal local

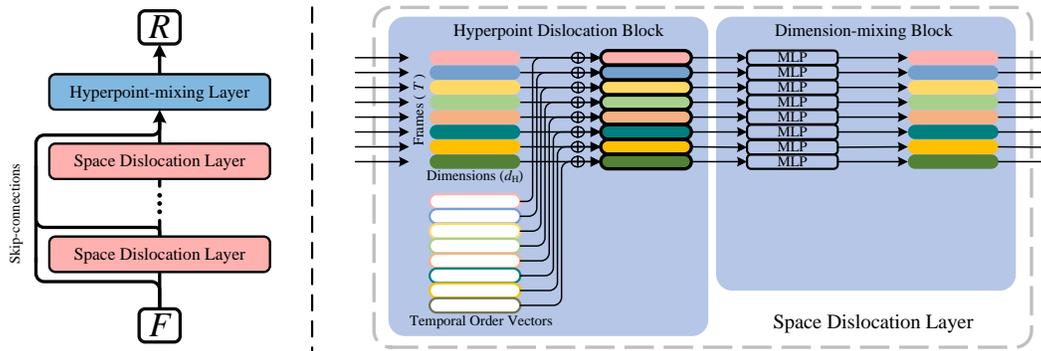


Fig. 4. Left: Hyperpoint-Mixer module. Right: Space dislocation layer.

structures to distinguish human actions. In our work, point cloud sequences are converted into Hyperpoint sequences to record the static appearances at each moment. To demonstrate the information equivalence property between point cloud sequences and Hyperpoint sequences, we provide a theoretical foundation for our converting operation by showing the universal approximation ability of the Hyperpoint embedding module to continuous functions on point cloud frames.

Formally, let $\mathcal{X} = \{S : S \subseteq [0, 1]^c \text{ and } |S| = n\}$ is the set of c -dimensional point clouds inside a c -dimensional unit cube. $f : \mathcal{X} \rightarrow \mathbb{R}$ is a continuous set function on \mathcal{X} w.r.t to Hausdorff distance $D_H(\cdot, \cdot)$, i.e., $\forall \epsilon > 0, \exists \delta > 0$, for any $S, S' \in \mathcal{X}$, if $D_H(\cdot, \cdot) < \delta$, then $|f(S) - f(S')| < \epsilon$. The theorem 1 [13] says that f can be arbitrarily approximated by PointNet given enough neurons at the max pooling layer.

Theorem 1. Suppose $f : \mathcal{X} \rightarrow \mathbb{R}$ is a continuous set function w.r.t Hausdorff distance $D_H(\cdot, \cdot)$. $\forall \epsilon > 0, \exists$ a continuous function h and a symmetric function $g(x_1, \dots, x_n) = \gamma \circ \text{MAX}$, such that for any $S \in \mathcal{X}$,

$$\left| f(S) - \gamma \left(\text{MAX}_{x_i \in S} \{h(x_i)\} \right) \right| < \epsilon$$

where x_1, \dots, x_n is the full list of elements in S ordered arbitrarily, γ is a continuous function, and MAX is a vector max operator that takes n vectors as input and returns a new vector of the element-wise maximum.

As stated above, continuous functions can be arbitrarily approximated by PointNet given enough neurons at the max pooling layer. The Hyperpoint embedding module is a recursive application of PointNet on nested partitions of the input point set. Therefore, the Hyperpoint embedding module is able to arbitrarily approximate continuous functions on point cloud frames given enough neurons at max pooling layers and a suitable partitioning strategy. In other words, the Hyperpoint embedding module is capable of extracting the complete static appearance information from point cloud frames to generate Hyperpoint sequences of equivalent information.

C. Hyperpoint-Mixer Module

Hyperpoint-Mixer module is proposed to model Hyperpoint sequence for 3D action recognition. The Hyperpoint sequence is an ordered set of high-dimensional single points. Compared to static point clouds and point cloud sequences,

modeling Hyperpoint sequences is extremely simple due to its ordered elements. In the Hyperpoint-Mixer module, we first dislocate Hyperpoints to specific temporal marking regions. Then, the dimension-mixing operation and the Hyperpoint-mixing operation are separately conducted to generate the final human spatio-temporal features. By doing so, the spatial and temporal information are decoupled, minimizing the impact of spatial irregularity on temporal information. Fig. 4 summarizes the Hyperpoint-Mixer module. The proposed Hyperpoint-Mixer module consists of multiple space dislocation layers of identical size, a multi-level feature learning based on skip-connection operations, and a Hyperpoint-mixing (\mathcal{H} -MIX) layer.

1) *Space dislocation layer*: The space dislocation layer is composed of a Hyperpoint dislocation block and a dimension-mixing (\mathcal{D} -MIX) block, which is presented to learn the internal spatial feature of dislocated Hyperpoints.

In the Hyperpoint dislocation block, Hyperpoints are spatially dislocated to record the temporal order by adding the corresponding temporal order vectors. The essence of the Hyperpoint dislocation block is to dislocate each Hyperpoint to a unique temporal marking region. In this fashion, the new coordinate of the dislocated Hyperpoint can be considered as the sum of two vectors: the temporal order vector and the spatial structure vector. The temporal order vector is used to distinguish the temporal order of the Hyperpoints while the spatial structure vector represents the internal appearance structure. Temporal order vectors (ToV) can be generated using the sine and cosine functions of different frequencies [15]:

$$ToV_{t,2h} = \sin \left(t/10000^{2h/d_H} \right) \quad (3)$$

$$ToV_{t,2h+1} = \cos \left(t/10000^{2h/d_H} \right) \quad (4)$$

where d_H denotes the dimension number of Hyperpoint coordinates. t is the temporal position and h is the dimension position.

Then, \mathcal{D} -MIX block learns the internal appearance structure of dislocated Hyperpoints in different temporal marking region, maps $\mathbb{R}^{d_H} \rightarrow \mathbb{R}^{d_H}$, and is shared across all dislocated Hyperpoints. Each \mathcal{D} -MIX block contains a set of MLP operations, a batch norm, and a ReLU non-linearity. Space

dislocation layers can be formalized as follows:

$$\mathbf{F}_t^\ell = \mathcal{D}\text{-MIX}(\mathbf{F}_t^{\ell-1} + \text{ToV}_t), \quad \text{for } t = 1 \dots T. \quad (5)$$

where \mathbf{F}^ℓ is the new Hyperpoint sequences after the ℓ -th space dislocation layer.

Each space dislocation layer takes an input of the same size, which is an isotropic design. With the stacking of space dislocation layers, Hyperpoints are dislocated at increasingly larger scales. In the larger dislocation space, the coordinates of Hyperpoints record more temporal information but less spatial information. To facilitate network optimization, multi-level features from different dislocation spaces are added by skip-connection operations and sent into \mathcal{H} -MIX layer as follows:

$$\mathbf{R}_i = \mathcal{H}\text{-MIX}_{t=1, \dots, T} \{(\mathbf{F} + \mathbf{F}^1 +, \dots, + \mathbf{F}^\ell)_{t,i}\} \quad (6)$$

where \mathbf{R} is the spatio-temporal feature of the Hyperpoint sequence.

2) *Hyperpoint-mixing layer*: \mathcal{H} -MIX layer is applied to aggregate spatial structures from all the Hyperpoints for the final spatio-temporal features. Since the temporal information has been injected, we can only use the simple max pooling to gather spatial information from the temporal marking regions. Benefiting from the order property of Hyperpoint sequences, temporal partitions can be easily divided. In order to capture the subactions within the Hyperpoint sequence, the hierarchical pyramid max pooling operation is adopted as the \mathcal{H} -MIX layer, which divides the fused Hyperpoint sequence $\tilde{\mathbf{F}} = \mathbf{F} + \mathbf{F}^1 +, \dots, + \mathbf{F}^\ell$ into multiple temporal partitions of the equal number of Hyperpoints and then performs the max pooling operation in each partition to generate corresponding spatio-temporal features. In this work, we employ a 2-layer pyramid with three partitions. Spatio-temporal features from all temporal partitions are simply concatenated to form the final spatio-temporal feature \mathbf{R} . Finally, the output of the Hyperpoint-Mixer module is input to a fully-connected classifier head for recognizing human actions.

IV. EXPERIMENTS

In this section, we firstly introduce the datasets. Then, we compare our SequentialPointNet with the existing state-of-the-art methods. Again, we conduct detailed ablation studies to further demonstrate the contributions of different components in our SequentialPointNet. After that, we compare the memory usage and computational efficiency of our SequentialPointNet with other point cloud sequence models. Finally, we apply our SequentialPointNet to skeleton sequence-based 3D human action recognition on NTU RGB+D 60 dataset.

A. Datasets

We evaluate the proposed method on two large-scale public datasets (*i.e.*, NTU RGB+D 60 [11] and NTU RGB+D 120 [16]) and a small-scale public dataset (*i.e.*, MSR Action3D dataset [10]).

The NTU RGB+D 60 dataset is composed of 56880 depth video sequences and skeleton video sequences for 60 actions

TABLE I
ACTION RECOGNITION ACCURACY (%) ON NTU RGB+D 60

Method/Year	Input	Cross-subject	Cross-view
Wang <i>et al.</i> (2018) [17]	depth	87.1	84.2
MVDI(2019) [18]	depth	84.6	87.3
3DFCNN(2020) [19]	depth	78.1	80.4
Stateful ConvLSTM(2020) [20]	depth	80.4	79.9
ST-GCN(2018) [21]	skeleton	81.5	88.3
AS-GCN(2019) [22]	skeleton	86.8	94.2
SGN(2020) [23]	skeleton	89	94.5
3s-CrosSCLR(2021) [24]	skeleton	86.2	92.5
Sym-GNN(2021) [25]	skeleton	90.1	96.4
3DV-PointNet++(2020) [5]	point	88.8	96.3
P4Transformer(2021) [8]	point	90.2	96.4
PSTNet(2021) [7]	point	90.5	96.5
HyperpointNet(2022) [9]	point	90.2	97.3
SequentialPointNet(ours)	point	90.3	97.6

TABLE II
ACTION RECOGNITION ACCURACY (%) ON NTU RGB+D 120

Method/Year	Input	Cross-subject	Cross-setup
Baseline(2018) [16]	depth	48.7	40.1
ST-GCN(2018) [21]	skeleton	81.5	88.3
MS-G3D Net (2020) [26]	skeleton	86.9	88.4
4s Shift-GCN(2020) [27]	skeleton	85.9	87.6
SGN(2020) [23]	skeleton	79.2	81.5
3s-CrosSCLR(2021) [24]	skeleton	80.5	80.4
3DV-PointNet++(2020) [5]	point	82.4	93.5
P4Transformer(2021) [8]	point	86.4	93.5
PSTNet(2021) [7]	point	87.0	93.5
HyperpointNet(2022) [9]	point	83.2	95.1
SequentialPointNet(ours)	point	83.5	95.4

and is one of the largest human action datasets. Both cross-subject and cross-view evaluation criteria are adopted for training and testing.

The NTU RGB+D 120 dataset is the largest dataset for 3D action recognition, which is an extension of the NTU RGB+D 60 dataset. The NTU RGB+D 120 dataset is composed of 114480 depth video sequences for 120 actions. Both cross-subject and cross-setup evaluation criteria are adopted for training and testing.

The MSR Action3D dataset contains 557 depth video samples of 20 actions from 10 subjects. Each action is performed 2 or 3 times by every subject. We adopt the same cross-subject settings in [10], where all the 20 actions are employed. Half of the subjects are used for training and the rest for testing.

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

In this section, in order to verify the recognition accuracy of our SequentialPointNet, comparison experiments with other state-of-the-art approaches are implemented on NTU RGB+D 60 dataset, NTU RGB+D 120 dataset, and MSR Action3D dataset.

1) *NTU RGB+D 60 dataset*: We first compare our SequentialPointNet with the state-of-the-art methods on the NTU RGB+D 60 dataset. As indicated in Table I, SequentialPointNet has recognition accuracy of 90.3% and 97.6% on the cross-subject and cross-view test settings, respectively. SequentialPointNet shows strong performance on par or even better than other point sequence-based approaches. Our SequentialPointNet achieves state-of-the-art performance among

TABLE III
ACTION RECOGNITION ACCURACY (%) ON MSR-ACTION3D

Method/Year	Input	# Frames	Accuracy
Kläser <i>et al.</i> (2008) [28]	depth	18	81.43
Vieira <i>et al.</i> (2012) [29]	depth	20	78.20
MeteorNet(2019) [6]	point	24	61.61
PointNet++(2020) [14]	point	1	88.50
P4Transformer(2021) [8]	point	24	90.94
PSTNet(2021) [7]	point	24	91.20
HyperpointNet(2022) [9]	point	24	91.54
SequentialPointNet(ours)	point	24	92.64

all methods on the cross-view test setting and results in similar recognition accuracy as PSTNet on the cross-subject test setting. The key success of our SequentialPointNet lies in the effective encoding for the temporal appearance evolution by Hyperpoint-Mixer module and minimizing the impact of the spatial irregularity on temporal information. Compared with our previous work HyperpointNet [9] that transforms Hyperpoint sequences into static point clouds, Hyperpoint-Mixer module in our SequentialPointNet designs isotropic space dislocation layers to blend the spatial and temporal information under multiple scales, enhancing the feature extraction ability for Hyperpoint sequences.

2) *NTU RGB+D 120 dataset*: We then compare our SequentialPointNet with the state-of-the-art methods on the NTU RGB+D 120 dataset. As indicated in Table II, SequentialPointNet achieves accuracy of 83.5% and 95.4% on the cross-subject and cross-setup test settings, respectively. Compared with PSTNet, SequentialPointNet does not show a competitive recognition accuracy on the cross-subject setting. However, SequentialPointNet gains a strong lead on the cross-setup test setting and achieves the highest recognition accuracy. Moreover, compared to existing point cloud sequence models, our approach greatly simplifies the point cloud sequence modeling task by decomposing it into a Hyperpoint embedding task and a Hyperpoint sequence encoding task.

3) *MSR Action3D dataset*: In order to comprehensively evaluate our method, comparative experiments are also carried out on the small-scale MSR Action3D dataset. To alleviate the overfitting problem on the small-scale dataset, the batch size is set as 16. Other parameter settings remain the same as those on the two large-scale datasets. Table III illustrates the recognition accuracy of different methods. When using 24 point cloud frames as inputs, our model yields state-of-the-art performance on the MSR Action3D dataset. Experimental results on the small-scale dataset demonstrate that our approach can achieve superior recognition accuracy even without a large amount of data for training.

C. Ablation Study

In this section, comprehensive ablation studies are performed on NTU RGB+D 60 dataset to validate the contributions of different components in our SequentialPointNet.

1) *Effectiveness of Hyperpoint-Mixer module*: We conduct the experiments to demonstrate the effectiveness of the Hyperpoint-Mixer module, and results are reported in Table IV. Several strong deep networks are used to instead of

TABLE IV
CROSS-VIEW RECOGNITION ACCURACY (%) WHEN USING MODELS

Method	Accuracy
SequentialPointNet(LSTM)	85.9
SequentialPointNet(GRU)	86.4
SequentialPointNet(Transformer)	81.6
SequentialPointNet(MLP-Mixer)	94.5
SequentialPointNet	97.6

TABLE V
CROSS-VIEW RECOGNITION ACCURACY (%) OF DIFFERENT METHODS

Method	Accuracy
SequentialPointNet(4D, w/o hdb)	95.4
SequentialPointNet(w/o mfl)	96.9
SequentialPointNet	97.6

the Hyperpoint-Mixer module in our SequentialPointNet. In SequentialPointNet (LSTM), LSTM is employed. In SequentialPointNet (GRU), GRU is employed. In SequentialPointNet (Transformer), a Transformer of two attention layers is used. In SequentialPointNet (MLP-Mixer), a MLP-Mixer of two mixer layers is used.

We can see from Table IV that results of SequentialPointNet (LSTM) and SequentialPointNet (GRU) are much worse when compared with SequentialPointNet. Hyperpoint sequences can be regarded as a kind of time series data. Different from the conventional time series data, the internal structures of elements rather than the element changes generate the main discriminant information. Thus, time series models that perform strict time-varying reasoning are not applicable to Hyperpoint sequences. Recently, self-attention-based Transformer has remained dominant in natural language processing [15] and computer vision [30]. However, due to the lack of larger-scale data for pre-training, SequentialPointNet (Transformer) does not show a competitive result. SequentialPointNet (MLP-Mixer) achieves the accuracy of 94.5%, which is 3.1% lower than our SequentialPointNet. The reason for this is that channel-mixing and token-mixing are performed alternatively, resulting in the mutual influence between spatial and temporal information.

2) *Different temporal information embedding manners*: The Hyperpoint dislocation block in the Hyperpoint-Mixer module is employed to embed the temporal information. To demonstrate its effectiveness, we also report the results of SequentialPointNet (4D, w/o hdb), which does not use the Hyperpoint dislocation block and injects the order information by appending the 1D temporal dimension to raw 3D points in each point cloud frame. Results are tabulated in Table V. From the table, we observe that SequentialPointNet with the Hyperpoint dislocation block outperforms SequentialPointNet (4D, w/o hdb). Therefore, the temporal information embedding manner used in SequentialPointNet is more efficient. From the experimental results, we can draw a conclusion that premature embedding of temporal information will affect spatial information encoding and decrease the recognition accuracy.

3) *Effectiveness of multi-level feature learning in the Hyperpoint-Mixer module*: With the stacking of the space dislocation layer in the Hyperpoint-Mixer module, Hyperpoints

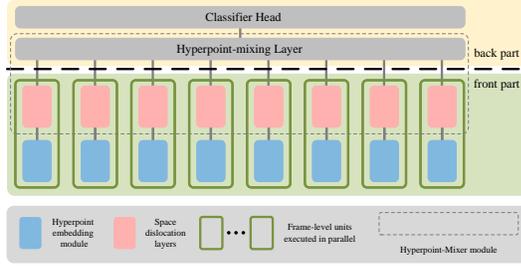


Fig. 5. The computation graph of SequentialPointNet.

are dislocated at increasingly larger scales. In the larger dislocation space, the coordinates of Hyperpoints record more temporal information but less spatial information. In order to obtain more discriminant information, multi-level features from different dislocation spaces are added by skip-connection operations. To verify the effectiveness of the multi-level feature learning, we report the result of SequentialPointNet (w/o mfl) that classify human actions without the multi-level features. We can see from Table V that the recognition accuracy of SequentialPointNet (w/o mfl) decreases by 0.7% without the multi-level feature learning.

D. Modeling Efficiency Analysis

In this section, We first show the frame-level parallelism of SequentialPointNet by its computation graph. Then, we evaluate the memory usage and computational efficiency of our method.

1) *Computation graph of SequentialPointNet*: Fig. 5 demonstrates the computation graph of SequentialPointNet. Taking the Hyperpoint-mixing layer as the watershed, SequentialPointNet can be divided into a front part and a back part. Since there is no cross-frame computational dependency, operations of the front part can be divided into frame-level units executed in parallel. Each frame-level unit includes a Hyperpoint embedding module and all space dislocation layers. The back part only contains the network architectures of low computational complexity including the Hyperpoint-mixing operation and a classifier head. Therefore, the main operations (*i.e.*, the front part) in SequentialPointNet can be executed in a frame-level parallel manner, based on which the modeling efficiency is greatly improved.

2) *Memory usage and computational efficiency*: Experiments are conducted on the machine with one Intel(R) Xeon(R) W-3175X CPU and one Nvidia RTX 3090 GPU. In Table VI, the number of parameters and the running time of our SequentialPointNet are compared with MeteorNet, 3DV-PointNet++, P4Transformer, and PSTNet. The running time is the network forward inference time per point cloud sequence.

From the table, we can see that parameters in 3DV-PointNet++ are the fewest in all methods. This is because 3DV-PointNet++ converts point cloud sequences to 3D point clouds, and employs static point cloud methods to process 3D point clouds. Compared with point cloud sequence methods, the static point cloud method has fewer parameters but lower recognition accuracy. The parameter number of SequentialPointNet is much less than MeteorNet, P4Transformer,

TABLE VI
PARAMETERS AND RUNNING TIMES COMPARISON

Method	Parameters (M)	Time (ms)
MeteorNet [6]	17.60	33.56
3DV-PointNet++ [5]	1.24	54.85
P4Transformer [8]	44.10	17.58
PSTNet [7]	8.44	27.56
SequentialPointNet	3.72	1.32

TABLE VII
COMPARISON BETWEEN SEQUENTIALPOINTNET AND ST-GCN ON NTU RGB+D 60

Method	Parameters (M)	Time (ms)	Accuracy (%)
ST-GCN [21]	3.12	4.62×10^{-1}	86.63
SequentialPointNet	0.55	1.84×10^{-2}	88.96

and PSTNet, which verifies that our approach is the most lightweight point cloud sequence model. In addition, SequentialPointNet is far faster than other methods. SequentialPointNet takes only 1.32 milliseconds to classify a point cloud sequence far beyond real-time requirement. In 3DV-PointNet++, the multi-stream network limits its parallelism. In MeteorNet, P4Transformer, and PSTNet, spatio-temporal local encoding are performed, which is time-consuming and not conducive to parallel computing. SequentialPointNet improves the speed of point cloud sequence modeling by more than 10 times. The superior computational efficiency of our SequentialPointNet is due to the lightweight network architecture and strong frame-level parallelism.

Additionally, our SequentialPointNet facilitates computational flexibility for point cloud sequence modeling. Since there is no computational dependency, in the case of limited computing power, frame-level units can also be deployed on different devices or executed sequentially on a single device.

E. SequentialPointNet for Skeleton Sequence-based 3D Human Action Recognition

SequentialPointNet can remarkably simplify the complexity of the sequence video data modeling task by decomposing it into a Hyperpoint embedding task based on corresponding existing static data encoding technology and a generic Hyperpoint sequence encoding task. To verify the effectiveness of SequentialPointNet's adaption to a skeleton sequence modeling task, we further apply our method to skeleton sequence-based 3D human action recognition on NTU RGB+D 60 dataset [11].

The frame in the skeleton sequence is a kind of graph structure data, representing the static appearance of the human body. Therefore, for the skeleton sequence, the existing graph convolution network (GCN) layers [21] is used to construct a spatial graph convolution network as the Hyperpoint embedding module of SequentialPointNet. The Hyperpoint-Mixer module remains the same as on the point cloud sequence. For a fair comparison, we sample 24 frames with equal intervals from a skeleton video sequence. We compare our SequentialPointNet with the spatial-temporal graph convolutional

network (ST-GCN) [21], which is the most famous GCN-based method for skeleton sequence-based action recognition. In Table VII, the number of parameters, the running time, and the cross-view accuracy of our SequentialPointNet are compared with ST-GCN. Skeleton sequences have a joint tracking nature that is convenient for exploring the joint-level movement information. Note that, SequentialPointNet without utilizing the joint tracking nature achieves even better recognition performance than ST-GCN. This is because that SequentialPointNet effectively encodes the temporally changing human appearances and minimizes the impact of spatial irregularity on temporal information. Furthermore, SequentialPointNet yields fewer network parameters and faster running time than ST-GCN.

V. CONCLUSION

In this paper, we propose a novel network named SequentialPointNet to model point cloud sequences for real-time 3D human action recognition. Instead of capturing spatio-temporal local structures, SequentialPointNet models the temporal evolution of static appearances to recognize human actions. In our SequentialPointNet, the point cloud sequence modeling is treated as a two-phase task: Hyperpoint embedding and Hyperpoint sequence modeling. We propose Hyperpoint as a new type of point cloud data and design a Hyperpoint-Mixer module as the basic building block to model Hyperpoint sequences for recognizing 3D human actions. Extensive experiments conducted on three public datasets show that SequentialPointNet obtains superior recognition performance and improves the speed of point cloud sequence classification by more than 10 times, significantly facilitating real-time 3D human action recognition based on point cloud sequences.

REFERENCES

- [1] Q. Zhu, Z. Chen, and C. S. Yeng, "A novel semi-supervised deep learning method for human activity recognition," *Transactions on Industrial Informatics*, vol. 17, no. 5, pp. 3821–3830, Jul. 2019.
- [2] T. Huynh-The, C. H. Hua, and D. S. Kim, "Encoding pose features to images with data augmentation for 3-D action recognition" *IEEE Transactions on Industrial Informatics*. 2019.
- [3] Z. Zuo, L. Yang, Y. Liu, F. Chao, R. Song, and Y. Qu, "Histogram of fuzzy local spatio-temporal descriptors for video action recognition," *IEEE Transactions on Industrial Informatics*. 2019.
- [4] Y. Min, Y. Zhang, X. Chai, and X. Chen, "An efficient pointlstm for point clouds based gesture recognition," *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 2020.
- [5] Y. Wang, Y. Xiao, F. Xiong, W. Jiang, Z. Cao, J. T. Zhou, and J. Yuan, "3dv: 3d dynamic voxel for action recognition in depth video," *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2020.
- [6] X. Liu, M. Yan, and J. Bohg, "MeteorNet: Deep Learning on Dynamic 3D Point Cloud Sequences," *Proceedings of the IEEE International Conference on Computer Vision*, 2019.
- [7] H. Fan, X. Yu, Y. Ding, Y. Yang, and M. Kankanhalli, "Pstnet: Point spatio-temporal convolution on point cloud sequences," *International Conference on Learning Representations*, 2021.
- [8] H. Fan, Y. Yang, and M. Kankanhalli, "Point 4d transformer networks for spatio-temporal modeling in point cloud videos," *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2021.
- [9] X. Li, Q. Huang, T. Yang, and Q. Wu, "HyperpointNet for point cloud sequence-based 3D human action recognition," *IEEE International Conference on Multimedia & Expo (ICME)*, 2022, doi: 10.1109/ICME52920.2022.9859807.
- [10] W. Li, Z. Zhang, and Z. Liu, "Action recognition based on a bag of 3d points," *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 2010.
- [11] A. Shahroudy, J. Liu, T.-T. Ng, and G. Wang, "Ntu rgb+d: A large scale dataset for 3d human activity analysis," *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 2016.
- [12] T. Le and Y. Duan, "PointGrid: A deep network for 3D shape understanding," *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 2018.
- [13] C. R. Qi, H. Su, K. Mo, and L. J. Guibas, "PointNet: Deep learning on point sets for 3D classification and segmentation," *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 2017.
- [14] C. R. Qi, L. Yi, H. Su, and L. J. Guibas, "Pointnet++: Deep hierarchical feature learning on point sets in a metric space," *Advances in Neural Information Processing Systems*, 2017.
- [15] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, and I. Polosukhin. "Attention is all you need," In *Advances in neural information processing systems*, 2017.
- [16] J. Liu, A. Shahroudy, M. Perez, G. Wang, L. Yudian, and A. C. Kot, "NTU RGB+D 120: A largescale benchmark for 3d human activity understanding," *IEEE transactions on pattern analysis and machine intelligence*, 2020.
- [17] P. Wang, W. Li, Z. Gao, C. Tang, and P. O. Ogunbona, "Depth pooling based large-scale 3-d action recognition with convolutional neural networks," *IEEE Transactions on Multimedia*, vol. 20, no. 5, 2018.
- [18] Y. Xiao, J. Chen, Y. Wang, Z. Cao, J. T. Zhou, and X. Bai, "Action recognition for depth video using multi-view dynamic images," *Information Sciences*, vol. 480, 2019.
- [19] A. Sanchez-Caballero, S. de López-Diz, D. Fuentes-Jimenez, C. Losada-Gutiérrez, M. Marrón-Romera, D. Casillas-Perez, and M. I. Sarker, "3dfcnn: Real-time action recognition using 3d deep neural networks with raw depth information," *arXiv preprint arXiv:2006.07743*, 2020.
- [20] A. Sanchez-Caballero, D. Fuentes-Jimenez, and C. Losada-Gutierrez, "Exploiting the convlstm: Human action recognition using raw depth video-based recurrent neural networks," *arXiv preprint arXiv:2006.07744*, 2020.
- [21] S. Yan, Y. Xiong, and D. Lin, "Spatial temporal graph convolutional networks for skeleton-based action recognition," *Thirty-second AAAI conference on artificial intelligence*, 2018.
- [22] M. Li, S. Chen, X. Chen, Y. Zhang, Y. Wang, and Q. Tian, "Actional-structural graph convolutional networks for skeleton-based action recognition," *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 2019.
- [23] P. Zhang, C. Lan, W. Zeng, J. Xing, J. Xue, and N. Zheng, "Semantics-guided neural networks for efficient skeleton-based human action recognition," *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2020.
- [24] L. Li, M. Wang, B. Ni, H. Wang, J. Yang, and W. Zhang, "3d human action representation learning via cross-view consistency pursuit," *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2021.
- [25] M. Li, S. Chen, X. Chen, Y. Zhang, Y. Wang, and Q. Tian, "Symbiotic graph neural networks for 3d skeleton-based human action recognition and motion prediction," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2021.
- [26] Z. Liu, H. Zhang, Z. Chen, Z. Wang, and W. Ouyang, "Disentangling and unifying graph convolutions for skeleton-based action recognition," *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2020.
- [27] K. Cheng, Y. Zhang, X. He, W. Chen, J. Cheng, and H. Lu, "Skeleton-based action recognition with shift graph convolutional network," *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2020.
- [28] A. Kläser, M. Marszałek, and C. Schmid, "A spatio-temporal descriptor based on 3d-gradients," *BMVC 2008-19th British Machine Vision Conference British Machine Vision Association*, 2008.
- [29] W. A. Vieira, E. R. Nascimento, G. L. Oliveira, Z. Liu, and M. F. Campos, "Stop: Space-time occupancy patterns for 3d action recognition from depth map sequences," *Iberoamerican congress on pattern recognition*, Springer, Berlin, Heidelberg, 2012.
- [30] A. Dosovitskiy, L. Beyer, and A. Kolesnikov, "An image is worth 16x16 words: Transformers for image recognition at scale," *arXiv preprint arXiv:2010.11929*, 2020.



Xing Li received the B.Sc. and M.Sc. degrees in software engineering from Changzhou University, China, in 2016 and 2019, respectively. He is currently pursuing the Ph.D. degree with the Department of Computer Science and Technology, Hohai University, Nanjing, China. His current research interests include machine learning, computer vision and deep learning, especially human action recognition.



Zhuang Miao received the Ph.D. degree from PLA University of Science and Technology, Nanjing, China. Zhuang Miao is currently a professor of Army Engineering University of PLA, Nanjing, China. His current research focuses on artificial intelligence, pattern recognition and computer vision.



Qian Huang received the B. Sc. degree in computer science from Nanjing University, China, in 2003, and the Ph. D. degree in computer science from the Institute of Computing Technology, Chinese Academy of Sciences, in 2010. From 2010 to 2012, he was a deputy technical manager of Mediatek (Beijing) Incorporation, Beijing, China. Since Dec. 2012, he is with Hohai University, Nanjing, China, where he serves as the dean of Computer Science & Technology Department. His research interests lie in industry-specific multimedia computing, especially on video compression & communication, object identification & tracking, and behavior expression & analysis. He is a member of the CAAI Technical Committee on Deep Learning, a member of the CCF Technical Committee on Multimedia Technology, and a member of the CSIG Technical Committee on Multimedia. Currently he serves as an associate editor for IET Image Processing, and a reviewer for some IEEE Transactions such as TIP, TMM and TCSVT.



Zhijian Wang received the Ph.D. degree from Nanjing University, NanJing, China. He is currently a Professor with Hohai University, Nanjing, China. His current research interests include machine learning, computer vision and deep learning.



Tianjin Yang received the B.Sc. and M.Sc. degrees in Computer Science and Technology from Changzhou University, China, in 2018 and 2021, respectively. He is currently pursuing the Ph.D. degree with the Department of Artificial Intelligence, Hohai University, Nanjing, China. His current research interests include machine learning, computer vision.



Zhenjie Hou received the Ph.D. degree in mechanical engineering from Inner Mongolia Agricultural University, in 2005. From 1998 to 2010, he was a Professor with the Computer Science Department, Inner Mongolia Agricultural University. In August 2010, he joined Changzhou University. His research interests include signal and image processing, pattern recognition, and computer vision.