

Learning Multiscale Correlations for Human Motion Prediction

Honghong Zhou¹, Caili Guo^{1,2*}, Hao Zhang¹ and Yanjun Wang³

¹Beijing Key Laboratory of Network System Architecture and Convergence,

School of Information and Communication Engineering,

Beijing University of Posts and Telecommunications, Beijing, China

²Beijing Laboratory of Advanced Information Networks, Beijing, China

³China Telecom Dict Application Capability Center, China

{zhouhonghong, guocaili, zhanghao215}@bupt.edu.cn, wangyanjun@chinatelecom.cn

Abstract—In spite of the great progress in human motion prediction, it is still a challenging task to predict those aperiodic and complicated motions. We believe that capturing the correlations among human body components is the key to understand the human motion. In this paper, we propose a novel multiscale graph convolution network (MGCN) to address this problem. Firstly, we design an adaptive multiscale interactional encoding module (MIEM) which is composed of two sub modules: scale transformation module (STM) and scale interaction module (SIM) to learn the human body correlations. Secondly, we apply a coarse-to-fine decoding strategy to decode the motions sequentially. We evaluate our approach on two standard benchmark datasets for human motion prediction: Human3.6M and CMU motion capture dataset. The experiments show that the proposed approach achieves the state-of-the-art performance for both short-term and long-term prediction especially in those complicated action category. We make codes publicly available at <https://github.com/zhouhongh/MGCN>.

Index Terms—Human motion prediction, multiscale, graph convolution network, DCT

I. INTRODUCTION

Human motion prediction aims to use the 3D skeleton data to predict a sequence of future human motions based on observed motion frames. It plays a significant role in robotics, computer graphics, healthcare and public safety [1]–[3] such as human robot interaction [4], autonomous driving [5], pedestrian tracking [6] etc.

Traditionally, Hidden Markov Model [7] and Gaussian Process latent variable models [8] is used to predict human motions, but limited to simple actions such as walking and golf swing. More complicated actions are typically tackled using deep networks including the recurrent neural networks (RNNs) [9]–[16] and feed-forward networks (FNNs) [17]–[24]. Simply apply RNN without the modeling of the body structure especially the correlations among human body suffered bad results [9]. Jain proposes to use the Structural RNN to model relationship among the spine and limbs which achieves good performance [15]. More and more people aware the importance of exploring the human body correlations from the body structure and proposed many RNN-based methods

[11]–[14]. Due to the RNN’s weakness on capturing long-term temporal dependencies and the disadvantage of error accumulation, the feed-forward networks have attracted more and more attentions. As in [20]–[23], the Discrete Cosine Transformation (DCT) becomes the popular temporal encoding strategy, which allows the network to concentrate on extracting the spatial correlations. Although achieved great progress, they all model the human motion in one single scale and the prediction on more complicated and aperiodic actions such as greeting and directing traffic is also a challenging task especially in the long-term scenario.

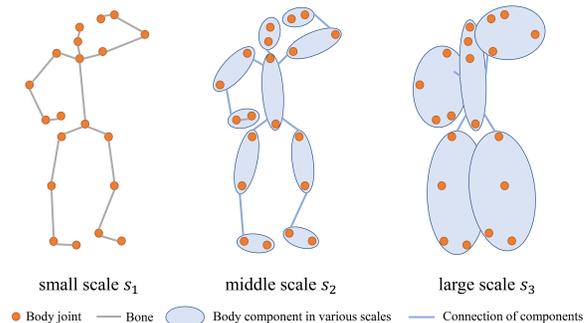


Fig. 1. Three body scales on Human 3.6M. In s_1 , we consider 20 joints with non-zero angles [11]; In s_2 and s_3 , we consider 10 and 5 components, respectively.

We keep under observation on motion rules of the real person, found that in some actions like running, the co-movement mainly exists among the big limbs, while in some other actions like smoking, the small movement of the wrist or elbow could lead to really different future poses. This scalable attention of human motions inspires us to capture the correlations of human body in a multiscale way. In this paper, based on the multiscale graph proposed in [24] shown in Fig. 1, we further exploit the multiscale modeling method and propose the multiscale graph convolution network (MGCN) as in Fig. 2, which achieve much better performance than [24].

In summary, our contributions are twofold:

- We make a comprehensive study on human motion prediction and propose an encoder-decoder framework called

*Corresponding author.

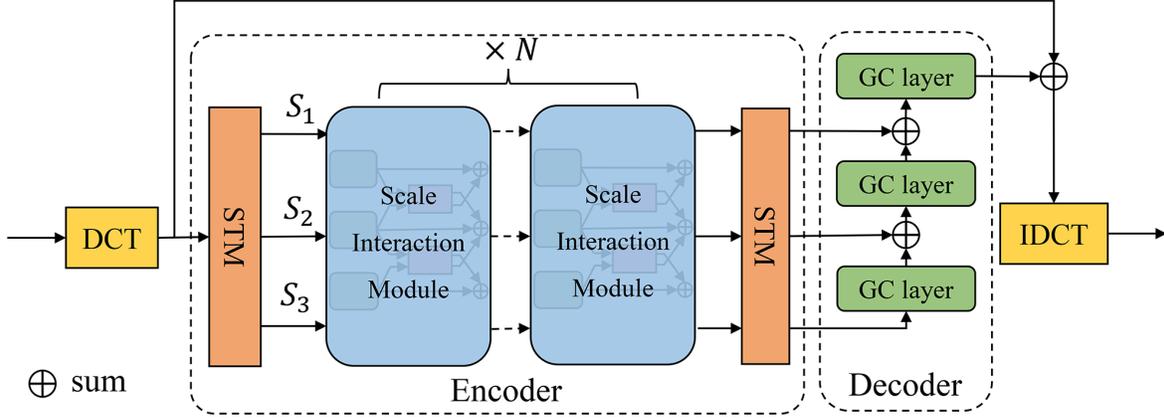


Fig. 2. The architecture of MGCN.

MGCN to deeply exploit the correlation of human body joints and components in the multiscale graphs, which is action-agnostic and end-to-end.

- We verify the effectiveness of MGCN on two benchmark datasets. Especially, our approach outperforms the SOTA method [22] 5 to 7 millimeters for the long-term prediction, which is a big lift.

II. RELATED WORK

In this section, we introduce the related works from the perspective of scales, including the single-scale and multiscale methods.

A. Single-scale methods

1) *RNN based methods*: The LSTM-3LR and ERD proposed by Fragkiadaki et al. [9] model the human motion based on concatenated LSTM units. S-RNN [15] models the main five components of human body by five RNNs and exchange information among them. After the S-RNN [15], more attention is paid on the capturing of the spatial dependencies especially the correlations of human body components. Liu et al. [11] design a hierarchical RNN which allows the information flows along adjacent joints and frames simultaneously. The evolution of these RNN based methods indicates the importance of modeling the human body correlations. While limited to RNN's drawbacks on error accumulation and long-term dependencies extracting, the FNN based methods become more popular these years.

2) *FNN based methods*: Li et al. [18] model the human motions with the feed-forward convolutional network and gain lower prediction errors than the existing RNN methods. Influenced by the interests in signal processing, Mao et al. [20] propose to learn the body joints trajectory dependencies with DCT and apply the fully-connected GCN to explore the correlations among body joints, which brings impressive progress. Inspired by Mao et al. [20], Lebailly et al. [21] introduce additional temporal inception module to achieve better temporal encoding. Cai et al. [21] replace the GCN

with the transformer network and achieves the state-of-the-art performance. The success of GCNs and transformer networks also shows the importance of modeling the human body correlations.

B. Multiscale methods

The multiscale modeling strategy is widely used in machine learning, such as object detection [27], [28], and NLP [29], [30]. It has a huge advantage over solving problems which have important features at multiple scales of time and/or space. On the human motion prediction task, Li et al. [24] propose to capture the human body correlations by the multiscale graph based on the backbone ST-GCN [31]. They generate human body graphs for three scales by mean-pooling, use the ST-GCN to encode information in each scale and design the cross-scale fusing blocks to fuse features with the adjacent scales. However, the mean-pooling strategy for generating multiscale graphs will cause information loss, and the ST-GCN mixes the temporal and spatial information, which is not benefit to the sequential human motion prediction task [32].

From the previous works we can draw a conclusion that the key to predict motions accurately is to properly model the human body correlations and the multiscale methods have great potential on human motion prediction. The mentioned background shows the big value of our work and we introduce it detailedly in the next section.

III. METHODOLOGY

On this task, we assume to be given a history motion sequence $\mathbf{X}_{1:N} = [\mathbf{x}_1, \mathbf{x}_2, \mathbf{x}_3, \dots, \mathbf{x}_N]$ consisting of N consecutive human poses, where $\mathbf{x}_i \in \mathbb{R}^K$, with K data dimensions describing pose at each time step. And our goal is to predict the future poses $\mathbf{X}_{N+1:N+T}$ for the future T time steps. Before sending the input data to the MIEM, following [20], we replicate the last pose \mathbf{x}_N , T times to generate a new sequence of length $N + T$: $\mathbf{X}'_{1:N+T} =$

$[\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_{N-1}, \mathbf{x}_N, \mathbf{x}_N, \dots, \mathbf{x}_N]$ and compute the DCT coefficients of length D for the new sequence as:

$$\mathbf{F} = f_{DCT}(\mathbf{X}'_{1:N+T}) \in \mathbb{R}^{K \times D} \quad (1)$$

where \mathbf{F} is the DCT coefficients, and f_{DCT} is the Discrete Cosine Transformation.

We make efforts to capture the human body correlations with the proposed MGCN whose architecture is shown in Fig. 2. We use those replicated DCT coefficients to predict the real ones and finally make IDCT to obtain the human motion frames on Euler angle representation or 3D coordinates.

A. Encoder

The MIEM plays the role of encoder. The MIEM includes two types of sub modules: 1) STM, which aggregates joints to components or convert components back into joints, and 2) SIM, which extract the human body correlations in and across three scales.

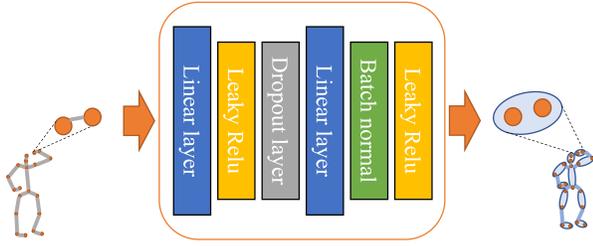


Fig. 3. STM core network. This figure shows how the STM transform 2 joints of s_1 into a component of s_2 . Actually it is a two-layer MLP, for each component of s_2 , we apply a MLP like this.

1) *Scale transformation module (STM)*: We design two STMs for our encoder according to Fig. 2, one is used to aggregate body joints into body components (before SIM) and the other is used to convert the aggregated body components back to body joints (after SIM). As shown in Fig. 3, the body joints at s_1 that belong to the same component are transformed to a new graph node at s_2 or s_3 . For example, we can aggregate the “right shoulder” joint and the “right elbow” joint to constitute the “right up arm” component. We apply a MLP for each component, since there are 10 body components at s_2 and 5 at s_3 , we apply totally 15 MLPs to realize the graphs’ scale transformation. Taking s_1 and s_2 as an example, the STM can be described by (2):

$$\mathbf{F}_k^{s_2} = f_1(\mathbf{F}_{i;j}^{s_1}) \quad (2)$$

where $\mathbf{F}_k^{s_2}$ indicates the features of the k^{th} node at s_2 , $\mathbf{F}_{i;j}^{s_1}$ is the features concatenated between the i^{th} and the j^{th} node at s_1 , f_1 is the two-layer MLP as Fig. 3. The mapping between k and i, j is defined in advance according to the human body structure.

When it comes to transform s_2 and s_3 back to the size of s_1 , the operation is similar. We just swap input dimension of MLPs with the output dimension. At s_2 , this process can be depicted by (3):

$$\mathbf{F}'_{i;j}{}^{s_2} = f_2(\mathbf{F}_k^{s_2}) \quad (3)$$

where $\mathbf{F}'_{i;j}{}^{s_2}$ is the features at s_2 but has the same number of node as that at s_1 . Noticed that $\mathbf{F}'_{i;j}{}^{s_2}$ still indicates the features of scale 2, we remain the s_2 superscript.

2) *Scale interactional module (SIM)*: In order to exploit the human body correlations more adequately, we design the SIM as Fig. 4, and cascade it N times to improve the feature extraction ability. the SIM is composed of two parts: the GCNs, which extract dynamic features in single scale, and the CS-Bs, which introduce the additional supervisory information from the adjacent scales. As we all know that if we want to recognize someone’s action, we need not only the cooperation among the large limbs such as arms and legs but also some subtle movements like the rotation of wrists, so the information-interacting strategy conforms our cognitive rules.

Following the notation of [20], we model the skeleton as a fully-connected graph of K nodes, represented by the trainable weighted adjacency matrix $\mathbf{A}^{K \times K}$. And the GCN is stacked by several GC layers, each performing the operation:

$$\mathbf{F}^{(p+1)} = \sigma(\mathbf{A}^{(p)} \mathbf{F}^{(p)} \mathbf{W}^{(p)}) \quad (4)$$

where $\mathbf{W}^{(p)}$ is the set of trainable weights of layer p , $\mathbf{A}^{(p)}$ is the learnable adjacency matrix of layer p , $\mathbf{F}^{(p)}$ indicates the input of layer p while $\mathbf{F}^{(p+1)}$ the input of layer $p+1$ (and the output of layer p), $\sigma(\cdot)$ is an activation function such as $\tanh(\cdot)$.

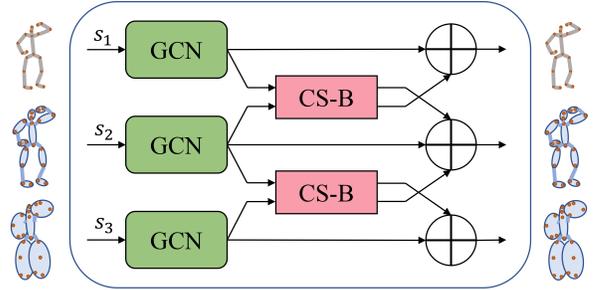


Fig. 4. The SIM network. The features of each scale get through a six-layer GCN, and relies on cross-scale blocks (CS-Bs) to explore the human correlations cross scales.

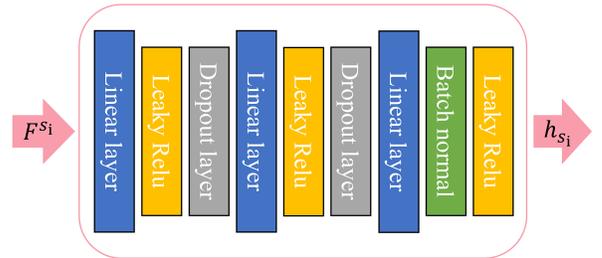


Fig. 5. The three-layer MLP that designed in the CS-B.

After the encoding of GCNs, the features interact cross scales by the cross-scale block (CS-B). Different from the

existing complicated operations in [24], we design three-layer MLPs as Fig. 5 to get attention matrix $\mathbf{A}_{s_i s_j}$ for two adjacent scales, which can speed up the calculation and be easy to train. The process to generate $\mathbf{A}_{s_i s_j}$ can be described by:

$$h_{s_1} = f_{MLP_1}(F^{s_1}) \in \mathbb{R}^{K_{s_1} \times D_h} \quad (5)$$

$$h_{s_2} = f_{MLP_2}(F^{s_2}) \in \mathbb{R}^{K_{s_2} \times D_h} \quad (6)$$

$$\mathbf{A}_{s_2 s_1} = \text{softmax}(h_{s_1}^T h_{s_2}) \in [0, 1]^{K_{s_2} \times K_{s_1}} \quad (7)$$

where f_{MLP_1} and f_{MLP_2} denote MLPs as Fig. 5, K_{s_i} means the number of nodes at s_i . $\mathbf{A}_{s_2 s_1}$ is the attention matrix from s_2 to s_1 . Benefit from the attention matrix $\mathbf{A}_{s_2 s_1}$, we can now adaptively explore the cross-scale human body correlations in a distinct way. We next introduce the supervisory information of the adjacent scale s_2 to s_1 with $\mathbf{A}_{s_2 s_1}$, the feature at s_1 is updated as:

$$\mathbf{F}^{s_1} \leftarrow \mathbf{A}_{s_2 s_1} \mathbf{F}^{s_2} + \mathbf{F}^{s_1} \in \mathbb{R}^{K_{s_1} \times D'} \quad (8)$$

B. Decoder

We apply the coarse-to-fine strategy to decode the features of three layers. The motivation of the decoder is that the larger scale can provide information of global motion evolution, which can indicate the approximate moving direction, speed and action category. The smaller scale then is able to predict the precise joint location with the global supervisory information. As in Fig. 2, the motion features at s_3 was sent to the bottom GC layer to predict the coarse future motion. Afterwards, we sum the output of the bottom GC layer and the feature at s_2 (middle scale) as the input of the middle GC layer. Similarly, the output of the top GC layer predicts the finest future motion. Mathematically, given the transformed features at s_2 and s_3 , \mathbf{F}'^{s_2} and \mathbf{F}'^{s_3} , and the feature at s_1 , \mathbf{F}^{s_1} , the predicted DCT coefficients \mathbf{F}_p are :

$$\mathbf{F}_p = f_{s_1}(f_{s_2}(f_{s_3}(\mathbf{F}'^{s_3}) + \mathbf{F}'^{s_2}) + \mathbf{F}^{s_1}) + \mathbf{F} \quad (9)$$

where f_{s_1} , f_{s_2} , f_{s_3} is the top, middle, bottom GC layer respectively as Fig. 4, \mathbf{F} is the original DCT coefficients.

Finally, we apply the IDCT to get the human motion frames on Euler angle representation or 3D coordinates.

C. Loss function

As in [20], when given the Euler angles, we use Mean Angle Error (MAE) as loss function to train our model, and use Mean Per Joint Position Error (MPJPE) proposed in [25] if given the 3D coordinates. Formally, the loss function on the angle and coordinate data can be described by (10) and (11), respectively:

$$\ell_a = \frac{1}{(N+T)K} \sum_{n=1}^{N+T} \sum_{k=1}^K |\hat{x}_{k,n} - x_{k,n}| \quad (10)$$

$$\ell_m = \frac{1}{J(N+T)} \sum_{n=1}^{N+T} \sum_{j=1}^J \|\hat{\mathbf{p}}_{j,n} - \mathbf{p}_{j,n}\|^2 \quad (11)$$

where $\hat{x}_{k,n}$ is the predicted k^{th} angles in frame n and $x_{k,n}$ the corresponding ground truth, $\hat{\mathbf{p}}_{j,n} \in \mathbb{R}^3$ denotes the predicted

j^{th} joint position at frame n , $\mathbf{p}_{j,n}$ is the corresponding ground truth, and J is the number of joints on the human skeleton graph.

IV. EXPERIMENTS

In this section, we introduce implementation details, followed by the datasets, the experimental results analysis and ablation study.

A. Implementation details

The GCN is the cascade of 6 residual blocks, each of which comprises 2 graph convolutional layers. We train the model for 100 epochs with a learning-rate decay of 0.96 every 2 epochs. The batch size on Human 3.6M dataset is 256 and On CMU dataset is 16. The stack number N of SIM is 3. The feature dim of GCN is set to 256. The feature dim of MLPs in the STM is 16 while that in the CS-B is 512. More implementation details can be found in our project home page at <https://github.com/zhouhongh/MGCN>.

B. Datasets

1) *Human3.6M*: There are 15 actions performed by 7 subjects for training and testing in the dataset. The actors are represented by a skeleton of 32 joints. Following the settings in [16], [26], we remove the global rotations and translations as well as constant angles and down sample the sequences to 25 frames per second.

2) *CMU-Mocap*: We select 8 actions and report results on the CMU mocap dataset (CMU-Mocap) following [18], [20], [22], [24]. It is also down sampled to 25 frames per second and removed the global rotations and translations as well as constant angles.

C. Baselines

We use 3 methods as the baselines: DMGNN [24], LTD [20] and LPJP [22]. The DMGNN applies the multiscale graph with the ST-GCN backbone. However, it does not adopt the timing modeling method of DCT, and does not use a strategy similar to STM and coarse-to-fine decoding. The LTD proposes the "DCT + space modeling" framework which uses the DCT to encode the trajectory of body joints and apply the fully-connected GCN to capture the human body correlations. But it completely ignores the multi-scale idea. The LPJP follows LTD's "DCT + spatial modeling" framework, and uses the popular Transformer [33] to replace the fully connected GCN and uses a central-to-peripheral ring decoding strategy to predict the human body, which makes LPJP the SOTA method.

Noted that the authors of DMGNN do not report their performance under the MPJPE metric and do not provide the model weight files, we train the DMGNN by ourselves following the settings in their paper and calculate the MPJPE. It is important to note that DMGNN does not contain codes to train directly on 3D coordinate data, we can only train on angle data and then calculate MPJPE. But this does not affect the comparison, because we also provide the MPJPE that calculated from our model trained on the angle data.

TABLE I
SHORT-TERM PREDICTION IN MAE ON HUMAN3.6M FOR THE MAIN ACTIONS.

milliseconds	Walking				Eating				Smoking				Directions				Greeting				Average			
	80	160	320	400	80	160	320	400	80	160	320	400	80	160	320	400	80	160	320	400	80	160	320	400
DMGNN [24]	0.18	0.31	0.49	0.58	0.17	0.30	0.49	0.59	0.21	0.39	0.81	0.77	0.25	0.44	0.65	0.71	0.36	0.61	0.94	1.12	0.27	0.52	0.83	0.95
LTD [20]	0.18	0.31	0.49	0.56	0.16	0.29	0.50	0.62	0.22	0.41	0.86	0.80	0.26	0.45	0.71	0.79	0.36	0.60	0.95	1.13	0.27	0.51	0.83	0.95
LPJP [22]	0.17	0.30	0.51	0.55	0.16	0.29	0.50	0.61	0.21	0.40	0.85	0.78	0.22	0.39	0.62	0.69	0.34	0.58	0.94	1.12	0.25	0.49	0.83	0.94
ours	0.18	0.30	0.47	0.51	0.16	0.28	0.46	0.57	0.21	0.38	0.82	0.76	0.23	0.40	0.73	0.80	0.34	0.57	0.91	1.08	0.26	0.49	0.81	0.92

TABLE II
SHORT-TERM PREDICTION IN MPJPE ON HUMAN3.6M FOR THE MAIN ACTIONS.

milliseconds	Walking				Eating				Smoking				Directions				Greeting				Average			
	80	160	320	400	80	160	320	400	80	160	320	400	80	160	320	400	80	160	320	400	80	160	320	400
DMGNN [24]	12.5	22.9	40.1	52.5	11.0	23.2	46.1	55.9	8.0	15.4	26.4	30.2	13.0	22.4	47.0	57.5	17.5	33.4	73.0	89.3	14.3	29.0	58.1	70.6
LTD [20]	11.1	19.0	32.0	39.1	9.2	19.5	40.3	48.9	9.2	16.6	26.1	29.0	11.2	23.2	52.7	64.1	14.2	27.7	67.1	82.9	13.5	27.0	54.0	65.0
LTD 3D [20]	8.9	15.7	29.2	33.4	8.8	18.9	39.4	47.2	7.8	14.9	25.3	28.7	12.6	24.4	48.2	58.4	14.5	30.5	74.2	89.0	12.1	25.0	51.0	61.3
LPJP [22]	9.6	18.0	33.1	39.1	9.1	19.5	40.2	48.8	7.2	14.2	24.7	29.7	9.3	22.0	51.6	63.2	15.4	30.7	71.8	82.8	11.9	26.1	53.2	64.5
LPJP 3D [22]	7.9	14.5	29.1	34.5	8.4	18.1	37.4	45.3	6.8	13.2	24.1	27.5	11.1	22.7	48.0	58.4	13.2	28.0	64.5	77.9	10.7	23.8	50.0	60.2
ours	10.2	18.0	30.5	36.6	8.9	18.7	37.5	46.0	7.5	13.6	23.1	27.4	9.5	20.1	48.5	58.3	14.4	28.5	62.7	77.6	12.5	25.5	50.8	61.8
ours 3D	8.1	15.0	27.1	31.3	8.2	18.4	37.7	44.5	6.8	12.9	22.6	27.0	10.6	22.7	48.9	60.1	12.8	25.8	68.4	86.8	10.8	23.2	49.3	60.1

TABLE III
LONG-TERM PREDICTION IN MPJPE ON HUMAN3.6M.

milliseconds	Walking		Eating		Smoking		Discussion		Average	
	560	1000	560	1000	560	1000	560	1000	560	1000
DMGNN [24]	56.5	90.1	80.6	107.7	39.6	61.2	96.4	111.2	68.3	92.6
LTD [20]	55.0	60.8	68.1	79.5	42.2	70.6	93.8	119.7	64.8	82.6
LTD 3D [20]	42.3	51.3	56.5	68.6	32.3	60.5	70.5	103.5	50.4	71.0
LPJP [22]	51.8	58.7	59.3	76.5	40.3	76.8	82.6	107.7	58.5	79.9
LPJP 3D [22]	36.8	41.2	58.4	67.9	29.2	58.3	74.0	103.1	49.6	67.6
ours	40.1	45.8	60.1	74.4	30.6	59.4	71.1	83.0	50.5	65.7
ours 3D	37.7	43.7	53.0	68.7	28.0	55.2	55.7	72.9	43.6	60.1

D. Results

For fair comparison, we report both short-term and long-term predictions for the two datasets, given the the history sequence including 10 frames as input. Noticed that the sequences have the speed of 25 frames per second, The short-term prediction means predicting for 400 milliseconds, 10 frames, and the long-term 1000 milliseconds, 25 frames.

1) *Results on Human 3.6M:* For short-term prediction, we evaluate our method under both MAE (Table I) and MPJPE (Table II) protocols, in comparison to state-of-the-art baselines [18], [20], [22], [24]. The 3D suffix to a method indicates the method is directly trained on 3D joint positions. Otherwise, the results were obtained by converting the joints angles to 3D positions. Due to the limited space, we report the results of the five main actions and the average results under all 15 actions. On the average results, we can see that our approach outperforms all the baselines at the later time steps, but a little worse than the SOTA method [22] at 80 milliseconds, which indicates that our approach mainly works at farther time steps. we speculate that it is because that at the closer time steps the movement is subtle and do not need too complicated spatial encoding. This problem also can be seen on walking and direction in Table II. In particular, let's look at the performance of DMGNN, which also uses a multi-scale strategy. Although DMGNN slightly outperformed other methods including ours method for a few moments, such as

320ms for Smoking in Table I, 320ms for Directions, and 400ms for Directions in Table II, our method leads the way in most categories, especially in the average metrics. It is proved that our method has significant advantages over DMGNN's multi-scale modeling methods, and even slightly superior to the current best method LPJP [22].

Additionally, our method is superior to LTD [20] in all categories in Table I, II and III, which shows that the single scale space modeling strategy is insufficient to capture the complicated human body correlations and it is necessary to introduce the multi-scale strategy in the "DCT + space modeling" framework. It needs to be pointed out that, as indicated in [20], the MAE metric may incorrectly evaluate the results on account of the cyclicity of angles, so we only report the results under MPJPE in the following experiments (long-term experiment on Human3.6M, short-term and long-term experiments on CMU and ablation study).

We also compare our results with baselines [20], [22], [24] in long-term scenarios in Table III. On the average metric, our method achieves a greater improvement than short-term scenarios and outperforms all baselines, indicating that our method can still maintain good accuracy in long-term prediction, thanks to our multi-scale strategy which fully captured human body correlations. However, in some actions, such as walking and eating, our method is slightly inferior to the SOTA method. We note that both walking and eating are strongly periodic movements, and such movements are easy to capture their body correlations. Therefore, the single-scale approach may be sufficient, while the improvement brought by the multi-scale strategy is not significant in these categories.

2) *Results on CMU-Mocap dataset:* To verify the universality of our approach, we also reported the results on CMU-Mocap dataset in both short and long-term scenarios by MPJPE as in Table IV. First, let's focus on the average indicator that best reflects the overall performance of the model. Looking at the methods with the 3D suffix, our method outperforms all baselines on average, which proves

TABLE IV
SHORT AND LONG-TERM PREDICTION IN MPJPE ON CMU-MOCAP DATASET.

	Basketball					Basketball Signal					Directing Traffic				
milliseconds	80	160	320	400	1000	80	160	320	400	1000	80	160	320	400	1000
DMGNN [24]	16.6	32.4	71.9	93.0	118.9	4.0	8.0	16.9	21.7	95.8	9.8	19.4	42.3	54.4	165.2
LTD 3D [20]	14.0	25.4	49.6	61.4	106.1	3.5	6.1	11.7	15.2	53.9	7.4	15.1	31.7	42.2	152.4
LPJP 3D [22]	11.6	21.7	44.4	57.3	90.9	2.6	4.9	12.7	18.7	75.8	6.2	12.7	29.1	39.6	149.1
ours	12.0	23.7	54.7	72.9	128.0	2.6	5.9	15.8	21.4	82.9	6.0	12.0	27.6	38.0	152.2
ours 3D	10.8	18.9	38.2	49.1	97.3	2.2	4.0	10.6	14.8	53.5	5.9	11.5	25.6	34.0	132.8
	Jumping					Running					Soccer				
milliseconds	80	160	320	400	1000	80	160	320	400	1000	80	160	320	400	1000
DMGNN [24]	21.4	44.3	96.0	118.7	191.2	11.3	17.87	21.7	26.3	69.5	15.4	32.3	68.0	80.7	167.9
LTD 3D [20]	16.9	34.4	76.3	96.8	164.6	25.5	36.7	39.3	39.9	58.2	11.3	21.5	44.2	55.8	117.5
LPJP 3D [22]	12.9	27.6	73.5	92.2	176.6	23.5	34.2	35.2	36.1	43.1	9.2	18.4	39.2	49.5	93.9
ours	13.9	29.8	77.9	102.4	177.7	21.2	29.6	27.2	28.7	75.7	9.1	18.2	42.6	53.3	125.2
ours 3D	13.4	29.5	74.0	96.9	162.1	17.4	21.2	20.6	26.5	65.1	9.1	16.7	37.5	52.5	119.5
	Walking					Washwindow					Average				
milliseconds	80	160	320	400	1000	80	160	320	400	1000	80	160	320	400	1000
DMGNN [24]	6.8	10.8	20.0	23.8	40.1	6.1	12.8	31.5	39.7	93.3	11.4	22.2	46.0	57.2	117.7
LTD 3D [20]	7.7	11.8	19.4	23.1	40.2	5.9	11.9	30.3	40.0	79.3	11.5	20.4	37.8	46.8	96.5
LPJP 3D [22]	6.7	10.7	21.7	27.5	37.4	5.4	11.3	29.2	39.6	79.1	9.8	17.6	35.7	45.1	93.2
ours	7.5	12.5	19.0	23.5	67.6	4.5	10.2	31.9	44.0	90.1	9.6	17.7	37.1	48.0	112.4
ours 3D	6.3	10.2	17.6	20.5	34.9	4.4	9.6	27.4	37.2	74.9	8.7	15.2	31.4	41.4	92.5

TABLE V
EFFECT OF THE MULTISCALE DESIGNS BY MPJPE ON HUMAN3.6M.

	Average			
milliseconds	80	160	320	400
w/o STM	13.1	26.8	52.9	63.5
w/o CS-B	12.6	25.8	51.9	63.1
w/o coarse-to-fine decoder	12.9	25.9	52.1	63.1
with all above	12.5	25.5	50.8	61.8

that the advantage of our method is stable and still has excellent performance on other data sets. And the results for methods without 3D suffix show that although both DMGNN and our approach adopt a multi-scale strategy, it is clear that our approach outperforms DMGNN on most categories, demonstrating that our multi-scale approach makes better use of the body associations of the human body.

In the long-term prediction of running, basketball, soccer, our method is somewhat inferior to the SOTA method, which may be because these categories all contain large periodic movements such as running, and can be well modeled without the need for multi-scale strategies. However, success in predicting more categories of actions and moments of human motions still shows significant advantages of our method.

E. Ablation study

We quantify the effect of our designs for the multiscale graph, including STM, CS-B and the coarse-to-fine decoder. The STM generates the larger graphs by adaptively aggregating the body joints to components. The CS-B exchange information among different scales, and the coarse-to-fine decoder decode the human motions sequentially. Noticed that directly

remove the STM will cause the collapse of the multiscale architecture, we replace it with the average strategy in [24] which simply forms the multiscale graphs by computing the mean value of body joints. Similarly, we replace the coarse-to-fine decoder by the simply parallel strategy which directly summing the outputs of three scales.

We train the model with all the 15 actions, and show the average value in Table V. We can see that the prediction errors increase no matter we remove any of the three parts. And the STM brings the biggest improvement, because it avoided the drop of information compared the average strategy in [24].

V. CONCLUSION AND FUTURE WORK

Human motion prediction has gained more and more attention with the rapid development of human-robot interaction and autonomous driving. Capturing of the human body correlations is the key to predict future motions. We propose the MGCN to explore the correlations by the multiscale graphs in and corss scales and our exhaustive experiments demonstrate that the proposed method outperform the state-of-the-arts methods especially for those complicated and aperiodic actions.

Further more, we note that current studies have focused on motion prediction for single person, with little consideration for multi-person scenarios. However, in real life, most human motions involve interactions with others, so it is obviously not sufficient to model a person without such interaction information. The next step of our work is to mine the interaction information among multiple persons, so as to achieve more accurate motion prediction in multi-person scenarios.

REFERENCES

- [1] Gupta, A., Martinez, J., Little, J.J., Woodham, R.J.: 3d pose from motion for crossview action recognition via non-linear circulant temporal encoding. In: Proceedings of the IEEE conference on computer vision and pattern recognition. pp. 2601-2608(2014)
- [2] Koppula, H., Saxena, A.: Learning spatio-temporal structure from rgb-d videos for human activity detection and anticipation. In: International conference on machine learning. pp. 792-800 (2013)
- [3] Kovar, L., Gleicher, M., Pighin, F.: Motion graphs. In: ACM SIGGRAPH 2008 classes. pp. 1-10 (2008)
- [4] Koppula, H.S., Saxena, A.: Anticipating human activities for reactive robotic response. In: IROS. p. 2071. Tokyo (2013)
- [5] Paden, B., C'ap, M., Yong, S.Z., Yershov, D., Frazzoli, E.: A survey of motion planning and control techniques for self-driving urban vehicles. *IEEE Transactions on intelligent vehicles* 1(1), 33-55 (2016)
- [6] Gong, H., Sim, J., Likhachev, M., Shi, J.: Multi-hypothesis motion planning for visual object tracking. In: 2011 International Conference on Computer Vision. pp. 619-626. IEEE (2011)
- [7] Brand, Matthew, and Aaron Hertzmann. "Style machines." Proceedings of the 27th annual conference on Computer graphics and interactive techniques. 2000.
- [8] Jack M Wang, David J Fleet, and Aaron Hertzmann. Gaussian process dynamical models for human motion. *IEEE transactions on pattern analysis and machine intelligence*, 30(2):283-298, 2008. 1
- [9] Fragkiadaki, K., Levine, S., Felsen, P., Malik, J.: Recurrent network models for human dynamics. In: Proceedings of the IEEE International Conference on Computer Vision. pp. 4346-4354 (2015)
- [10] Wang, B., Adeli, E., Chiu, H.k., Huang, D.A., Niebles, J.C.: Imitation learning for human pose prediction. In: Proceedings of the IEEE International Conference on Computer Vision. pp. 7124-7133 (2019)
- [11] Z. Liu et al., "Towards Natural and Accurate Future Motion Prediction of Humans and Animals," 2019 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR), Long Beach, CA, USA, 2019, pp. 9996-10004, doi: 10.1109/CVPR.2019.01024.
- [12] Hu J, Fan Z, Liao J, et al. Predicting Long-Term Skeletal Motions by a Spatio-Temporal Hierarchical Recurrent Network[J]. arXiv preprint arXiv:1911.02404, 2019.
- [13] Tang Y, Ma L, Liu W, et al. Long-term human motion prediction by modeling motion context and enhancing motion dynamic[J]. arXiv preprint arXiv:1805.02513, 2018.
- [14] Pavlo D, Grangier D, Auli M. Quaternet: A quaternion-based recurrent model for human motion[J]. arXiv preprint arXiv:1805.06485, 2018.
- [15] Jain, Ashesh, et al. "Structural-rnn: Deep learning on spatio-temporal graphs." Proceedings of the IEEE conference on computer vision and pattern recognition. 2016.
- [16] Liang-Yan Gui, Yu-Xiong Wang, Xiaodan Liang, and Jose MF Moura. Adversarial geometry-aware human motion prediction. In ECCV, pages 786-803, 2018. 1, 2, 5, 6, 7
- [17] Judith Butepage, Michael J. Black, Danica Kragic, and Hedvig Kjellstrom. Deep representation learning for human motion prediction and classification. In CVPR, July 2017. 1, 2, 3
- [18] Chen Li, Zhen Zhang, Wee Sun Lee, and Gim Hee Lee. Convolutional sequence to sequence model for human dynamics. In CVPR, pages 5226-5234, 2018. 1, 2, 5, 6, 7
- [19] Pang, Bo, et al. "Complex sequential understanding through the awareness of spatial and temporal concepts." *Nature Machine Intelligence* 2.5 (2020): 245-253.
- [20] Mao, Wei, et al. "Learning trajectory dependencies for human motion prediction." Proceedings of the IEEE/CVF International Conference on Computer Vision. 2019.
- [21] Lebailly, Tim, et al. "Motion Prediction Using Temporal Inception Module." Proceedings of the Asian Conference on Computer Vision. 2020.
- [22] Cai, Yujun, et al. "Learning progressive joint propagation for human motion prediction." *European Conference on Computer Vision*. Springer, Cham, 2020.
- [23] Mao, Wei, Miaomiao Liu, and Mathieu Salzmann. "History repeats itself: Human motion prediction via motion attention." *European Conference on Computer Vision*. Springer, Cham, 2020.
- [24] Li, Maosen, et al. "Dynamic multiscale graph neural networks for 3d skeleton based human motion prediction." Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. 2020.
- [25] Catalin Ionescu, Dragos Papava, Vlad Olaru, and Cristian Sminchisescu. Human3.6m: Large scale datasets and predictive methods for 3d human sensing in natural environments. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 36(7):1325-1339, jul 2014. 4, 5, 6
- [26] Julieta Martinez, Michael J. Black, and Javier Romero. On human motion prediction using recurrent neural networks. In CVPR, July 2017. 1, 2, 4, 5, 6, 7
- [27] Li, Yanghao, et al. "Scale-aware trident networks for object detection." Proceedings of the IEEE/CVF International Conference on Computer Vision. 2019.
- [28] Zhao, Hengshuang, et al. "Pyramid scene parsing network." Proceedings of the IEEE conference on computer vision and pattern recognition. 2017.
- [29] Vig, Jesse. "A Multiscale Visualization of Attention in the Transformer Model." Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics: System Demonstrations. 2019.
- [30] White, Daniel A., et al. "Multiscale topology optimization using neural network surrogate models." *Computer Methods in Applied Mechanics and Engineering* 346 (2019): 1118-1135.
- [31] Yan, Sijie, Yuanjun Xiong, and Dahua Lin. "Spatial temporal graph convolutional networks for skeleton-based action recognition." Proceedings of the AAAI conference on artificial intelligence. Vol. 32. No. 1. 2018.
- [32] Pang, Bo, et al. "Complex sequential understanding through the awareness of spatial and temporal concepts." *Nature Machine Intelligence* 2.5 (2020): 245-253.
- [33] Vaswani A, Shazeer N, Parmar N, et al. Attention is All you Need[C]/NIPS. 2017.