

Temporal Consistency Two-Stream CNN for Human Motion Prediction

Jin Tang*, Jin Zhang and Jianqin Yin*

School of Artificial Intelligence, Beijing University of Posts and Telecommunications, Beijing, China.

ARTICLE INFO

Keywords:

Temporal fusion
Two-stream network
Human motion Prediction

ABSTRACT

Fusion is critical for a two-stream network. In this paper, we propose a novel temporal fusion (TF) module to fuse the two-stream joints' information to predict human motion, including a temporal concatenation and a reinforcement trajectory spatial-temporal (TST) block, specifically designed to keep prediction temporal consistency. In particular, the temporal concatenation keeps the temporal consistency of preliminary predictions from two streams. Meanwhile, the TST block improves the spatial-temporal feature coupling. However, the TF module can increase the temporal continuities between the first predicted pose and the given poses and between each predicted pose. The fusion is based on a two-stream network that consists of a dynamic velocity stream (V-Stream) and a static position stream (P-Stream) because we found that the joints' velocity information improves the short-term prediction, while the joints' position information is better at long-term prediction, and they are complementary in motion prediction. Finally, our approach achieves impressive results on three benchmark datasets, including H3.6M, CMU-Mocap, and 3DPW in both short-term and long-term predictions, confirming its effectiveness and efficiency.

1. Introduction

Human motion prediction has been a classic task in the field of computer vision and robotics. It demands predicting the future motion postures of the human body by observing the previous motion sequence, which helps the robots judge the intention of humans more accurately and achieve more excellent human-machine interaction applications [1, 2, 3]. Unlike motion recognition [4, 5], which only demands modeling the semantic information of the human body, the challenge of human motion prediction is modeling the dynamics of the human body.

Most of the earliest traditional methods that deal with the human motion prediction task have adopted hidden Markov models [6] and linear dynamic systems [7], and so on. They work well in "Walking", "Eating" and other regular movements with high repetition patterns. However, it is challenging to learn the pattern of complex actions such as "Walking Dog", "Posing", and so on. After the emergence of deep learning methods [8, 9], models can represent high dimensional information, like features from different convolution layer extraction or derivative operated inputs. How to efficiently use and fuse these high-dimensional features has become a significant issue.

For enriching the dynamic temporal information, the optical flow information is used as dynamic temporal information at the pixel level in the video recognition. Sarma et al. [10] modeled the optical flow as an additional stream to model the target's motion characteristics in the image. And Wang et al. [11] considered the human motion velocity as one of three inputs to feed into the network to model the human body dynamics for motion prediction. Therefore, we believe that the velocity may have advantages in modeling

the joint-level motion dynamics.

Two-stream network [15, 16] has become a common network structure to introduce additional information. In this structure, the features are extracted from the two streams separately and then fused together to get output so that the information contained in the fusion features is more prosperous than that in every single stream. Therefore, constructing an effective fusion has become the main problem to be solved in the two-stream network.

For feature fusing, many works [12, 13, 14, 15, 16] propose to apply multi-scale features or high-dimensional information to model human dynamics. Some [12, 13, 15, 16] use feature addition to add features from a residual connection. Besides, Li et al. [14] used convolution layers to fuse the concatenation output of multi-scale features. However, as shown in Fig. 1(a)(b), methods like concatenation or addition neglect spatial-temporal co-occurrence and continuity of time. In this paper, we design a novel temporal fusion module that can ensure the coupling of the two-stream information by restoring the two channels of features' time dimension, as shown in Fig. 1(c).

As noted above, we model the velocity and position prediction from two different streams respectively. And the future velocity sequence of human joints generated by the V-stream will be restored into the position space to maintain its physical significance equal to that of position space. Then, we conduct the TF module for fusing position space output and velocity space output so as to take advantage of the dynamic information of human motion and give the ultimate prediction results. The temporal concatenation fuses V-stream features with P-stream features separately by the dimension of time, which achieves better integration of spatiotemporal information comparing with feature addition. In this way, each time-step's P-stream features and V-stream features are fused separately to ensure temporal consistency. To solve the problem of poor coupling after fusion, we use a TST module to continue modeling the spatiotemporal cou-

This work was supported partly by the Fundamental Research Funds for the Central Universities(Grant No. 2020XD-A04-1).

*Corresponding author

✉ tangjin@bupt.edu.cn (J. Tang); jinzhang@bupt.edu.cn (J. Zhang);
jqyin@bupt.edu.cn (J. Yin)

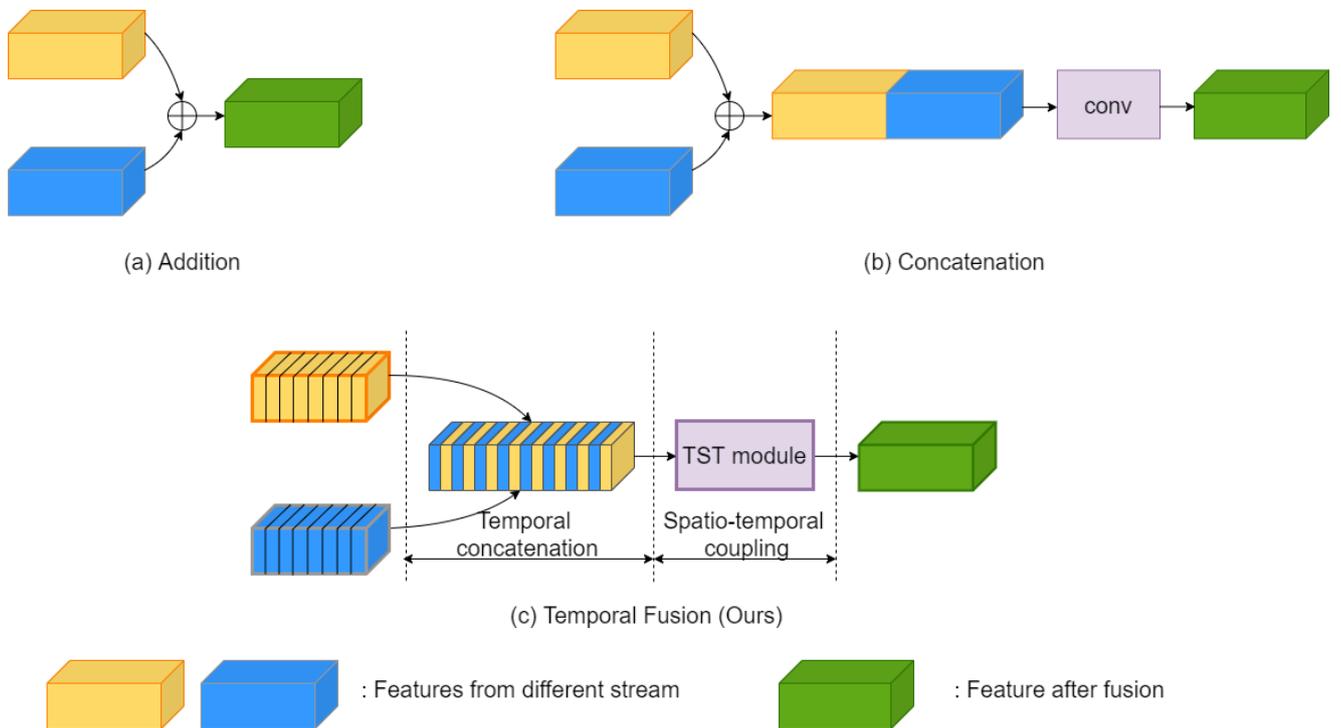


Fig.1: Methods for feature fusion. (a) Addition of two-stream features. (b) Concatenation of two-stream features. (c) A temporal fusion module based on temporal concatenation and a reinforcement trajectory spatial-temporal (TST) module. When inputting features from different modeling streams, (a)(b) neglect the temporal consistency comparing with (c).

pling features and generate the ultimate prediction results.

In summary, the main contributions are two-fold: (1) The novel temporal fusion (TF) module that keeps the temporal consistency and enhances the spatial-temporal features coupling. (2) A two-stream framework that integrates the advantages of velocity stream in short-term modeling and position stream in long-term modeling. Our approach is also general and easy to incorporate into a two-stream-based prediction framework. Our experiments on three standard human motion prediction benchmarks evidence the benefits of our approach.

The remainder of the paper is organized as follows. The next section investigates the related work. Section 3 discusses our model in detail. The dataset, evaluation criteria, the experimental-based comparisons of different methods, and ablation studies are presented in Section 4. Finally, conclusions and future work directions are stated in section 5. Our code will be made public after the paper is accepted.

2. Related work

Since human motion can be regarded as a temporal sequence, recurrent neural network (RNN) methods [12, 17, 18, 19, 20] have achieved quite a remarkable effect on human motion prediction. Gopalakrishnan et al. [17] proposed a two-stage processing RNN architecture to capture the difference in the power spectrum of the ground truth frames, which helped reduce an accumulation of next-step error generated recursively. Along with Fragkiadaki et al. [19] pro-

posed to use Encoder-Recurrent-Decoder (ERD) to extend previous long-short-term memory (LSTM) models in the literature to jointly learn representations and their dynamics. However, [17, 18, 19] neglect the problem of discontinuities between the observed poses and the predicted future ones. Martinez et al. [13] used a sequence-to-sequence residual RNN to solve discontinuities. However, it still neglected the co-occurrence of the spatial-temporal relationship from a human motion sequence. The recursive modeling ability of the RNN network is excellent, but most of the works that performing RNN ignore the correlation between the spatial information of human joints and time. Liu et al. [20] proposed hierarchical motion context modeling to update the local skeletal state and introduced the spatial connection between joint points into an RNN. However, the proposed network loses sight of the global temporal co-occurrence relationship of the motion sequence.

Instead of RNN, convolutional neural networks (CNN) achieve better performance in many works [14, 21, 22]. They adopt a sequence to sequence model that performs convolution operation in the temporal domain so that the global temporal feature of human motion correlation can be captured. Li et al [14] used a long-term convolutional encoder to capture the long-term information. Spatial-temporal features are fed into the next layer at a different level to model the motion dynamic effectively. Liu et al. [21] also proposed a high-efficiency spatial-temporal modeling network, noted as TrajectoryNet. TrajectoryNet is a feed-forward network,

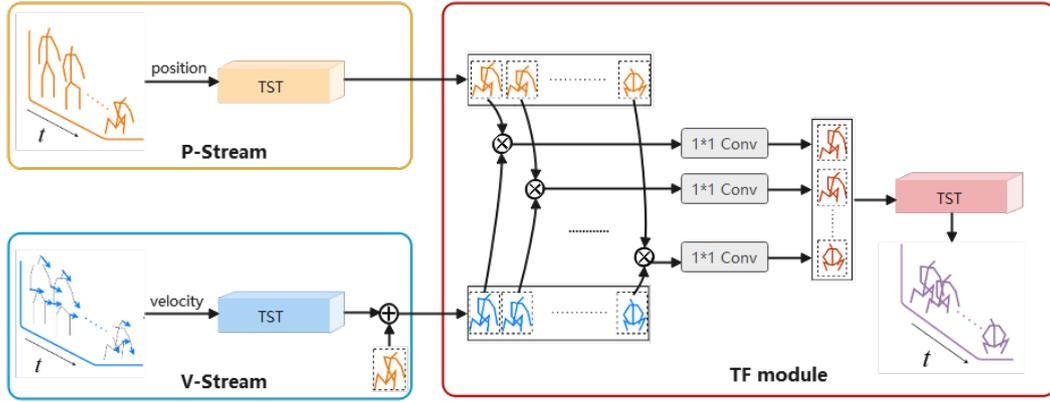


Fig. 2: Overall Architecture. The model of velocity is also important for skeleton-based motion prediction but is neglected in most earlier works. In this work, a V-Stream and a P-Stream are used to model velocity and position respectively. The TF module is composed of temporal concatenation and a reinforcement TST block.

treating the human motion as a composition of joint-level trajectories and modeling them by designing joints in a natural skeletal order, which retained both temporal and spatial features. Both works in [14, 21] model human dynamics in position space, conducting 3D joint-level coordinates as input and leaving high dimensional space input to be discussed.

Optical flow [23, 24, 25, 26] is widely used high-dimensional information for the task of video action recognition. It searches the corresponding relationship between the previous frame and the current frame by the difference of pixels between adjacent frames, which can be denoted as the pixel-wise velocity. As for joint-wise, recently, many works [11, 12, 13, 18, 19, 20, 27, 28] choose to model human motion dynamics by multi-level dynamic information to predict the regular movement changes of human motion and achieved well-performance prediction. Gui et al. [18] proposed an adversarial geometry-aware encoder-decoder (AGED) to perform adversarial training at the sequence level, which adopts the idea of multi-level GAN. The discriminators designed in AGED are used to discriminate the fidelity and continuity of predicted sequence separately. And Gui et al. [27] proposed to adopt proactive and adaptive meta-learning (PAML) to jointly learns a generic model initialization and an effective model adaptation strategy. [11] directly points out that the velocity of the input sequence can be presented as input to better model human dynamism by an RNN network. Human pose, pose velocity and position embedding are fed into RNN composed by GRU to encode different frames' absolute temporal positions and preserve motion continuities. Shu et al. [28] proposed a spatiotemporal co-attention RNN to model dynamics in skeleton motion and joint motion corresponding to temporal evolution and spatial coherence. In summary, multi-level dynamic information can be explored to enrich learning features for dif-

ferent tasks.

In many works [10, 11, 15, 16], two-stream or multi-stream modeling networks are proposed to effectively integrate high-dimensional or multi-level information into position space features and achieve better recognition or prediction results. Shi et al. [16] proposed a two-stream GCN-based framework to model both the bone and joint information respectively to increase the flexibility of the model for graph construction and bring more generality to adapt to various data samples. As well as Wang et al. [11] designed a Position-Velocity RNN (PVRNN) that took in three inputs and predicts pose velocities, which are then added to the previous posts to get the future poses. The spatial-temporal co-occurrence of the learned feature is obviously ignored. Then the predictions of pose velocities and human poses are used as input for the next time step in the RNN modeling process. Furthermore, [10, 15, 16] also lose sight of the part of fusion. The human motion sequence features are simply fused by linear options. Inspired by this, this paper, based on CNN for spatial-temporal modeling, proposes a fusion module suitable for human motion temporal sequence, to fuse the features from position space and velocity space generated by a two-stream network.

3. Methodology

The human motion prediction task is to observe the input pose sequence, which is denoted as $S = \{P_1, P_2, \dots, P_{t_i}\}$ along the dimension of time, and generate future pose sequence generated which is denoted as $S' = \{P_{t_i+1}, P_{t_i+2}, \dots, P_{t_i+t_o}\}$. The i th pose P_i is composed of joints, which can be denoted as $P_i = \{J_1, J_2, \dots, J_n\}$. Index n means the number of human skeletal joints. Meanwhile, the k th joint can be represented as $J_k = \{x_k, y_k, z_k\}$ in 3D coordinate space.

However, for the specific task of human motion predicting, there are a few common pitfalls that we would like to improve like the predicted pose tend to converge to the mean pose or fail to generate natural-looking poses due to clear discontinuities in the first frame. In this paper, we attack the above problems in a variety of ways. First, we introduce the velocity feature vector, computed by a time-step position difference, feeds into predictive neural models, which naturally holds local temporal information that is crucial when generating smooth and consistent motion trajectories. We predict the human pose by adding the predicted velocity on the last observed frame. We observe that naive input replacement cannot achieve high performance in long-term prediction even though it effectively improves the short-term prediction and gets more natural-looking poses. Furthermore, we then propose a two-stream convolutional network that consisting of a dynamic velocity stream and a static position stream, as shown in Fig. 2. To keep the spatial-temporal coherence, we propose a TF module to fuse prediction pose from P-Stream and V-Stream by the dimension of time, which achieves higher predictive performance in both short-term and long-term predictions. The details of the architecture are discussed below.

3.1. Architecture

We use a convolutional network to build the prediction pipeline since convolutional network can better model the spatial-temporal coupling features. The framework of our network consists of two parts, as shown in Fig. 2.

(1) For the first part, we model position space information and velocity space information in two streams using the TST block, respectively. Then each stream gives preliminary predictions. The predicted velocity vector was accumulated to the last frame of position space input to recover the velocity output to position space for fusion.

(2) For the second part, we fuse the prediction results from two streams in chronological order by our temporal concatenation. Since simple proportional fusion by 1*1 convolution still leads to poor spatial-temporal coupling, we use another TST block to continue modeling the spatial-temporal co-occurrence features.

In summary, the TST block is a vital block in our network. We use it in the two-stream modeling phase and TF module. The illustration of the TST is shown in Fig. 3.

Trajectory Spatial-temporal(TST). The TST block aims to model the trajectory of each input tensor. Since human motion can be regarded as composed of each joint’s trajectory, the convolution layer can model its trajectory information by the time dimension as the channel. Firstly, we process H3.6M as 22 key joints in the human body for training which is similar to the settings as the baseline TrajectoryNet [21]. Then, According to VGG [29], two 3*3 convolutions have the same receptive field as one 5*5 convolution when stride and padding are set to 1. We employ eleven 3*3 convolution layers in the TST block. The final equivalent receptive field of temporal dimension is enlarged, capable of covering

the size of a skeletal pose. Meanwhile, we use six 1*1 convolutional layers to construct residual connections, as shown in Fig. 3. The deeper layer is connected with the lower layer by a residual connection so that the coarse-grained features can be supplied with fine-grained features.

3.2. Dynamic Velocity Stream (V-Stream).

Many works have proposed methods of using velocity to model the dynamics of objects. As mentioned in Introduction, optical flow describes the displacement of pixels, which can be regarded as pixel-level velocity. Karen et al. [23] used a two-stream network to process video RGB frames and optical flow respectively to give the classification confidence. Since video is a temporal composition of RGB frames, human motion can also be treated as a temporal composition of 3D joint coordinates. Inspired by this, the concept of velocity is applied to represent the joint-level velocity information in our work.

Since the observed data is the joints’ position, the motion sequence’s velocity should be extracted from the observed position sequence. We adopt the difference between adjacent time-steps of a joint’s position to indicate the joint’s velocity at the moment. That is to say, the velocity of the k th joint at the moment t $V_{(k,t)}$ can be calculated by Eq [1]:

$$V_{(k,t)} = \{x_{k,t+1} - x_{k,t}, y_{k,t+1} - y_{k,t}, z_{k,t+1} - z_{k,t}\} \quad (1)$$

In this case, $V_{k,t}$ captures more dynamic temporal information and discards redundant information from position space at the moment of t . The input sequence of velocity can be described as $S_v\{V_1, V_2, \dots, V_{t-1}\}$. As shown in Fig.2, we denote the velocity sequence’s shape as $(N \times T - 1 \times 3)$ to represent the velocity of each joint of the human body at different moments, and the position sequence’s shape as $(N \times T \times 3)$. N is the number of joints and T is the number of input frames, "3" contains 3D coordinates of x, y, z . This constructed high-dimensional information will be piped into the V-stream network to generate a prediction that pays more attention to dynamic information. The prediction of velocity can be represented as $\widehat{S}_v = \{\widehat{U}_v^1, \widehat{U}_v^2, \dots, \widehat{U}_v^i\}$, and the predicted velocity of the joints in the i th frame \widehat{U}_v^i is composed by $\{\widehat{V}_{(1,i)}, \widehat{V}_{(2,i)}, \dots, \widehat{V}_{(n,i)}\}$, and the j th joint’s velocity at the t th frame can be represented as $\widehat{V}_{(j,t)} = \{\widehat{\Delta x}, \widehat{\Delta y}, \widehat{\Delta z}\}$. These deviation vectors are superposed to the last frame of the position sequence to recover it into position space for further fusion. That is, we calculate the i th predicted pose of the velocity stream by Eq [2].

$$\widehat{P}_v^i = \{J_{1,t} + \sum_{k=1}^i \widehat{V}_{1,t}, J_{2,t} + \sum_{k=1}^i \widehat{V}_{2,t}, \dots, J_{n,t} + \sum_{k=1}^i \widehat{V}_{n,t}\} \quad (2)$$

This preliminary prediction from the velocity stream will be used for feature fusing in the next section.

3.3. Temporal Fusion module (TF module)

Since simple concatenation ignores the temporal correlation between the position space and the velocity space, the

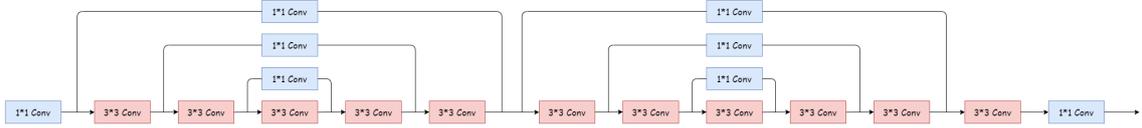


Fig.3: Trajectory spatial-temporal (TST) block. It improves the spatial-temporal feature coupling by using residual connections. Each convolution layer is coupling with leaky-Relu and Dropout option.

TF module comprises a temporal concatenation and a reinforcement trajectory spatial-temporal (TST) network. For temporal concatenation, a typical approach is applied, which is recovering the predicted sequence length \hat{T} from the number of channels, then they are fused in chronological order. As shown in Fig. 2, each stream outputs a tensor at the shape of $(N \times \hat{T} \times 3)$ as its prediction before the fusion. Since the tensor could be restored to the time dimension, the fusion based on each moment can be conducted respectively. The prediction results generated from P-stream and V-stream will be respectively divided into \hat{T} tensors in chronological order, and then the same time-step tensors are concatenated. In this way, we achieve \hat{T} tensors at the size of $(N \times 2 \times 3)$. For each moment's temporal concatenation result, one 1×1 convolution layer as a dynamic selector to balance the prediction from two streams. That is to say, these \hat{T} 1×1 convolution layers give a preliminary prediction fusing result of position space and velocity space for each predicting moment. In this way, each selector keeps its weight for leveraging two-stream features. Comparing with many works discussed in [10, 11, 15, 16] that fuse two-stream features by concatenation or addition, our method maintains the temporal correlation between the position space and the velocity space.

As shown in Fig.2, the temporal concatenation result is sent to a TST block in the TF module. Since the 1×1 convolution layer only gives the proportional fusion result, there is a lack of spatial coupling that each joint uses the same proportion to calculate the new joint position. In case of this, an additional TST block is added to refine the spatial-temporal coupling and output the final fusion result.

4. Experiments

We evaluate our model and compare it with the state-of-the-art on three benchmark datasets, including Human3.6M(H3.6M) [30], the CMU mocap dataset (CMU-Mocap), and the 3D pose in the Wild dataset (3DPW) [31]. In the following text, the experiment details, datasets, the evaluation metrics, and the baseline are first introduced, then the results of our method are presented and analyzed, and the predictive performance of our model is also visualized in the last.

4.1. Datasets

H3.6M. H3.6M [30] dataset is a commonly used dataset for human motion prediction, captured and recorded by 7 professional actors for training and evaluation. 15 activities are provided by pose sequences. Each frame is represented as 32 joints in relative 3D joint positions. Following existing works [14, 21], we down-sample the sequences by 2 to 25 frames per second and use Subject 5 as the test data and Subjects 1, 6, 7, 8, 9 as training. Preprocessing and data selection criteria are directly followed from the recent work of [14].

CMU-Mocap. CMU-Mocap¹ provides 2,235 human motion sequences. The frames are recorded by VGA cameras. Following baselines in [12, 21], eight actions are selected from category "locomotion", "physical activities & sports", "common behaviors and expressions" and "communication gestures and signals", we use the same dataset splits for training and testing.

3DPW. 3DPW [31] consists of indoor and outdoor actions such as shopping, doing sports, and hugging, including 60 sequences and more than 51k frames. For a fair comparison, we use the official split sets for experiments.

4.2. Implementation Details

Following the setting and processing in [21], all experiments are carried out in 3D coordinate space. In experiments, all models are implemented by TensorFlow. To be consistent with the literature, we report our results for short-term (< 400 ms) and long-term (> 400 ms) predictions. We follow the same experiment setting in [12, 14, 21], giving 10 frames (400 milliseconds) as input to predict 10 or 25 frames for short-term or long-term evaluation on H3.6M and CMU-Mocap. For 3DPW, we use 30 frames (1 second) as long-term prediction. Meanwhile, we use the same skeletal representation, dropout option and activation function as TrajectoryNet [21] for modeling local spatial features. All models are trained with Adam optimizer, and the learning rate is initialized to 0.0001. MPJPE (Mean Per Joints Position Error) [30] is used as our loss function to train the network, see Eq [3].

$$l = \frac{1}{T_o * N} \sum_{t=1}^{T_o} \sum_{k=1}^N \|\hat{J}_{t,k} - J_{t,k}\|^2 \quad (3)$$

¹<http://mocap.cs.cmu.edu/>

Table 1: Short-term prediction of 3D coordinates on H3.6M for all actions. Our method outperforms the baselines on average prediction on all time steps.

Milliseconds	Walking				Eating				Smoking				Discussion			
	80	160	320	400	80	160	320	400	80	160	320	400	80	160	320	400
ConvSeq2Seq[14]	17.1	31.2	53.8	61.5	13.7	25.9	52.5	63.3	11.1	21.0	33.4	38.3	18.9	39.3	67.7	75.7
LearnTrajDep[12]	8.9	15.7	29.2	33.4	8.8	18.9	39.4	47.2	7.8	14.9	25.3	28.7	9.8	22.1	39.6	44.1
TrajectoryNet[21]	8.2	14.9	30.0	35.4	8.5	18.4	37.0	44.8	6.3	12.8	23.7	27.8	7.5	20.0	41.3	47.8
Ours	7.8	14.8	28.9	34.4	7.7	17.1	34.5	42.6	6.2	12.7	23.3	27.6	7.1	18.7	39.7	46.8
Milliseconds	Directions				Greeting				Phoning				Posing			
	80	160	320	400	80	160	320	400	80	160	320	400	80	160	320	400
ConvSeq2Seq[14]	22.0	37.2	59.6	73.4	24.5	46.2	90.0	103.1	17.2	29.7	53.4	61.3	16.1	35.6	86.2	105.6
LearnTrajDep[12]	12.6	24.4	48.2	58.4	14.5	30.5	74.2	89.0	11.5	20.2	37.9	43.2	9.4	23.9	66.2	82.9
TrajectoryNet[21]	9.7	22.3	50.2	61.7	12.6	28.1	67.3	80.1	10.7	18.8	37.0	43.1	6.9	21.3	62.9	78.8
Ours	9.4	22.9	53.7	65.1	12.7	28.0	64.6	79.1	10.0	18.6	37.9	44.4	6.8	21.6	63.5	79.9
Milliseconds	Purchases				Sitting				Sitting Down				Taking Photo			
	80	160	320	400	80	160	320	400	80	160	320	400	80	160	320	400
ConvSeq2Seq[14]	29.4	54.9	82.2	93.0	19.8	42.4	77.0	88.4	17.1	34.9	66.3	77.0	14.0	27.2	53.8	66.2
LearnTrajDep[12]	19.6	38.5	64.4	72.2	10.7	24.6	50.6	62.0	11.4	27.6	56.4	67.6	6.8	15.2	38.2	49.6
TrajectoryNet[21]	17.1	36.1	64.3	75.1	9.0	22.0	49.4	62.6	10.7	28.8	55.1	62.9	5.4	13.4	36.2	47.0
Ours	17.4	36.1	59.4	67.4	8.6	21.5	50.3	63.6	11.0	27.4	51.2	60.7	5.5	13.5	37.1	49.0
Milliseconds	Waiting				Walking Dog				Walking Together				Average			
	80	160	320	400	80	160	320	400	80	160	320	400	80	160	320	400
ConvSeq2Seq[14]	17.9	36.5	74.9	90.7	40.6	74.7	116.6	138.7	15.0	29.9	54.3	65.8	19.6	37.8	68.1	80.2
LearnTrajDep[12]	9.5	22.0	57.5	73.9	32.2	58.0	102.2	122.7	8.9	18.4	35.3	44.3	12.1	25.0	51.0	61.3
TrajectoryNet[21]	8.2	21.0	53.4	68.9	23.6	52.0	98.1	116.9	8.5	18.5	33.9	43.4	10.2	23.2	49.3	59.7
Ours	8.3	20.5	51.4	66.3	20.8	47.7	94.1	108.7	7.7	18.1	31.9	41.1	9.8	22.6	48.1	58.4

$J_{t,k}$ denotes the k th joint’s 3D coordinates at the t th moment from the ground truth. Similarly, $\hat{J}_{t,k}$ denotes the k th joint’s 3D coordinates at the t th moment from prediction results. T_o and N represents the length of the output sequence and the number of joints. MPJPE can measure the average value of Euclidean space errors between the ground truth and the predicted joints and is often used to evaluate human motion estimation and prediction results. We do not adopt MPJPE on the preliminary prediction results generated by the two streams independently because we intend to avoid an excessive degree of feature convergence learned by the two streams. We notice that the regularity of 3DPW outdoor actions is poor, then we set the weight of the last 22 predicted frames to 0.2 for reducing error accumulating while training.

We also make use of the Mean Per Joint Position Error(MPJPE) [30] in millimeters as the error evaluation metric on all datasets. As mentioned in [14], angles are not a good representation to evaluate motion prediction. We employ the measurement of the Euclidean distance between the ground-truth pose and our predicted pose in the 3D coordinate space as the error metric. For example, the H3.6M dataset contains several human motion pose sequences. After dividing the training set and the test set, we take $T_i + T_o$ frames of the sequence as a sample. The first T_i frames are input to our two-stream network to obtain the predicted sequence \hat{S} which is composed by predicted joint $\hat{J}(t, k)$, and

the last T_o frames of the sample consist the ground truth sequence S . Then the MPJPE between the \hat{S} and the ground truth sequence S is calculated according to Eq[3].

4.3. Baselines

We compare our method with three convolutional networks, LearnTrajDep [12], convSeq2Seq [14] and TrajectoryNet [21]. The MPJPE results are taken from their respective papers. Specifically, LearnTrajDep [12] uses a graph convolutional network to learn graph connectivity automatically. ConvSeq2Seq [14] proposes a convolutional encoder-decoder network to capture both invariant and dynamic information of human motion. The proposed encoder extracts multi-level spatial-temporal features for forwarding modeling. TrajectoryNet [21] proposes to learn human dynamics in trajectory space by CNN modeling and residual connection. We implement the same data preprocessing with the baselines, which takes 3D coordinates as input and output.

4.4. Results and Discussion

This section presents the experimental results on the datasets, and a related discuss is provided.

Results on H3.6M. In Table 1, we compare our model with baselines for short-term prediction. Note that our two-stream model outperforms all the baselines on average prediction errors, demonstrating the effectiveness of our

Table 2: Long-term prediction of 3D coordinates on H3.6M for all actions. Our method outperforms the baselines on average prediction on all time steps.

Milliseconds	Walking		Eating		Smoking		Discussion		Directions		Greeting		Phoning		Posing		
	560	1000	560	1000	560	1000	560	1000	560	1000	560	1000	560	1000	560	1000	
ConvSeq2Seq[14]	59.2	71.3	66.5	85.4	42.0	67.9	84.1	116.9									
LearnTrajDep[12]	42.2	51.6	57.1	69.5	32.5	60.7	70.5	99.6	79.6	102.9	95.8	89.9	62.6	113.8	107.2	211.8	
TrajectoryNet[21]	37.9	46.4	59.2	71.5	32.7	58.7	75.4	103	84.7	104.2	91.4	84.3	62.3	113.5	111.6	210.9	
Ours	37.7	48.3	56.3	71.7	29.6	58.1	78.6	104.6	80.3	97.3	87.6	85.0	58.7	112.0	113.9	213.9	
Milliseconds	Purchases		Sitting		Sitting Down		Taking Photo		Waiting		Walking Dog		Walking Together		Average		
560	1000	560	1000	560	1000	560	1000	560	1000	560	1000	560	1000	560	1000	560	1000
ConvSeq2Seq[14]																	
LearnTrajDep[12]	92.4	125.3	78.6	109.9	88.1	137.8	78.3	95	99.0	169.5	139.2	167.7	60.3	84.1	78.9	112.6	
TrajectoryNet[21]	84.5	115.5	81.0	116.3	79.8	123.8	73.0	86.6	92.9	165.9	141.1	181.3	57.6	77.3	77.7	110.6	
Ours	86.0	114.5	81.4	116.2	82.4	126.1	69.8	89.8	91.2	162.9	135.4	166.7	55.3	77.6	76.3	109.6	

proposed approach. Comparing with the results of convSeq2Seq [14] and LearnTrajDep [12], our model predicts much better on movement actions. LearnTrajDep [12] respectively model the temporal information and the spatial dependencies of joint trajectories using GCNs, neglecting the co-relation of spatial-temporal features. However, our approach reinforces the spatial-temporal coupling by trajectory spatial-temporal (TST) block, which models joint-level trajectory. Therefore, our model achieves higher accuracy using convolutional TST module along with V-stream dynamic features. Also, comparing with the results of TrajectoryNet [21], our method achieves lower error on some specific common actions such as "Walking", "Eating" and "Smoking". These possible reasons are two folds: (1) The predicted velocity information accumulated to the last frame of the input sequence can keep the motion continuity. (2) Our temporal fusion (TF) module harmoniously keeps the temporal consistency of the P-stream and V-stream feature and better mine the law of the motion due to the fusion of the two streams. So, the accuracy of prediction results combined with two-stream features is higher, especially for regular moves like "Walking" or big move like "Walking Dog".

We provide qualitative comparisons in Fig. 4. Our model effectively captures human motion dynamics and shows better prediction results on "Walking", "Walking Dog". However, the baselines show discontinuity at the first predicting frame, and our method is more natural. These movements are inseparable from the regular swing of the legs so that the velocity vector can retain the displacement of the legs and the velocity stream can model the dynamics of legs more effectively. As for the other motions such as "Direction", "Taking Photo", our model doesn't give a better prediction than TrajectoryNet. To analyze our approach's failure cases, we visualize the poses in Fig. 5. As for the action "Taking Photo", there are irregular arm movements in these actions that may disturb velocity modeling. Consequently, the irregular arm motion pattern learned from the observation sequence harms the prediction results. However, as shown in "Smoking" for comparison, our model can still give proper prediction in some situ actions.

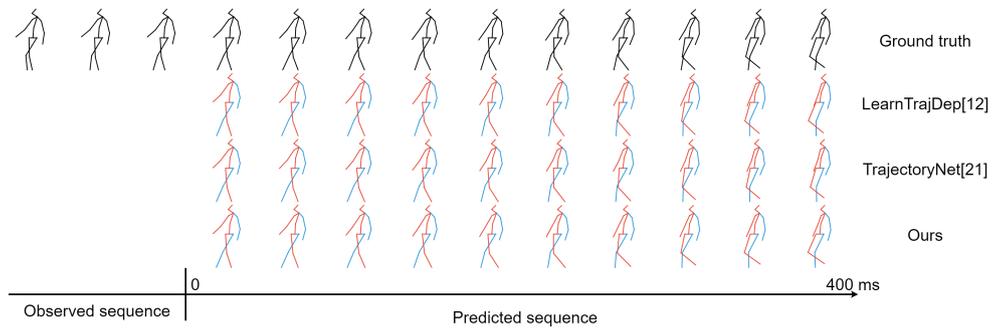
Our long-term prediction result on H3.6M is shown in Table 2. In long-term experiments, the prediction length is set as 25 frames corresponding to the duration of 1 second. Our model still achieves lower average errors comparing with the baselines. Especially for the movements that perform much better in our model on short-term prediction like "Walking Dog", the long-term prediction still achieves lower error. Therefore, we believe our TF module benefits human dynamics modeling.

Results on 3DPW. The results of short-term and long-term prediction on 3DPW dataset are shown in table 3. Our method shows better predicting accuracy on short-term predictions. However, the prediction results at 800ms and 1000ms do not exceed the TrajectoryNet [21]. 3DPW dataset is a wild action dataset, unlike H3.6M, which has special categories for regular actions such as "Walking" and "Eating", the regularity of 3DPW outdoor actions is complex, so it is difficult for our velocity stream to model the joint movement pattern of the action sequence, and it will accumulate more errors in long-term prediction results.

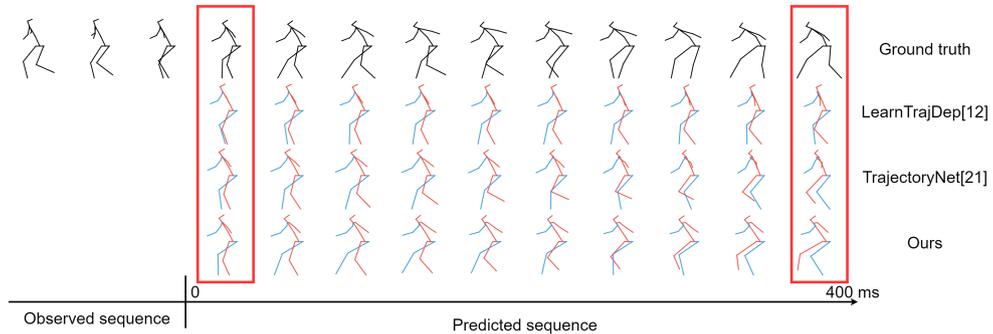
Results on CMU-Mocap. Our short-term and long-term prediction result on CMU-Mocap is reported in Table 4. Our method outperforms the baselines on both short-term and long-term predictions. The errors decrease slightly on all time-steps.

Discussion. As shown in Table 4, our model outperforms the baselines on all time-steps average predictions in H3.6M and CMU-Mocap and underperforms the TrajectoryNet long-term prediction in 3DPW. The reasons are analyzed for two aspects as below.

(1) Dataset analysis. H3.6M and CMU-Mocap record the positional data of human skeleton joints collected by professional equipment indoors. The center of the movement is located at the waist of the human body. Our proposed method can powerfully capture the dynamics of the movement actions and achieve state-of-the-art predictive performance in these datasets. Besides, H3.6M contains more movement actions than CMU-Mocap, which explains our

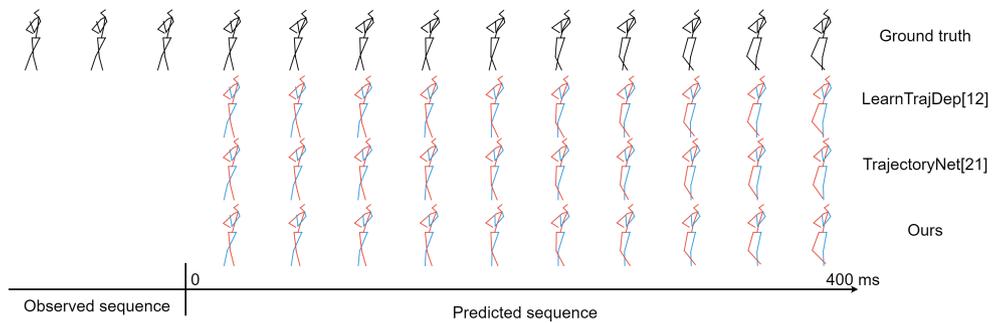


(a) Walking

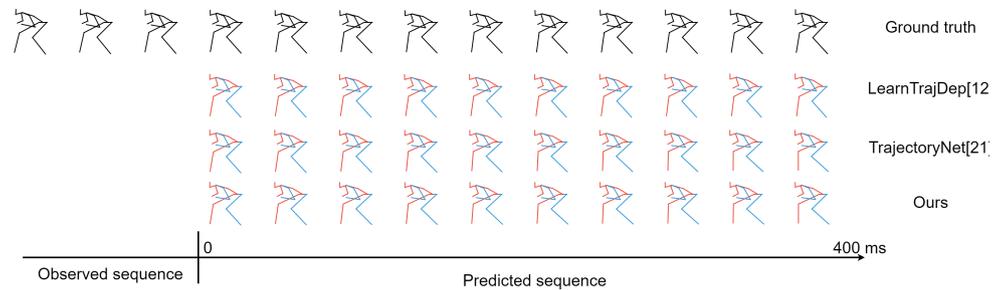


(b) Walking Dog

Fig.4: Quality comparison of "Walking" and "Walking Dog". Both are movement actions. Top: the conditioning sequence and the ground-truth of the predicted sequence. Middle two: state-of-the-art prediction results. Bottom: our prediction. The ground-truth and the input sequences are shown in black. The results evidence that our approach generates high-quality predictions in both cases.



(a) Smoking



(b) Taking Photo

Fig.5: Quality comparison of "Smoking" and "Taking Photo". Both are in situ actions. Top: the conditioning sequence and the ground-truth of the predicted sequence. Middle two: state-of-the-art prediction results. Bottom: our prediction. The ground-truth and the input sequences are shown in black.

Table 3: Short-term and long-term prediction on 3DPW

Milliseconds	200	400	600	800	1000
ConvSeq2Seq[14]	71.6	124.9	155.4	174.7	187.5
LearnTrajDep[12]	35.6	67.8	90.6	106.9	117.8
TrajectoryNet[21]	30.0	59.7	85.3	99.0	107.7
Ours	28.4	59.1	84.6	100.0	108.2

Table 4: Short-term and long-term prediction on CMU-Mocap

Milliseconds	80	160	320	400	1000
LearnTrajDep[12]	11.5	20.4	37.8	46.8	96.5
TrajectoryNet[21]	8.3	15.6	33.4	43.1	92.8
Ours	8.2	15.1	32.8	43.0	92.6

superior results on H3.6M. 3DPW is an outdoor recorded human action dataset. Therefore, the movements of 3DPW are more disturbed and less regular, which makes the prediction more challenging. As mentioned in results on 3DPW, because of the poor movement regularity of 3DPW’s long-term action samples, the velocity features interfere with the position space features, so the long-term prediction results do not exceed TrajectoryNet.

(2) Method analysis. Our proposed two-stream model introduces velocity as an auxiliary input to model human dynamics. The velocity vector can easily capture the regular human motion pattern, which explains our short-term prediction advantage. For further prediction, our model mostly depends on the learned regular pattern which consists of observed velocity features. Therefore, when coupling some irregular movements shot in outdoor locations, the velocity features lead to error accumulating. That explains our superiority in movement actions and disadvantage on some situ actions and irregular actions. After constructing velocity vectors, we use the TF module to fuse the two-stream features chronologically to ensure temporal consistency. In contrast, TrajectoryNet only uses spatial information for modeling, and its dynamic information is not as rich as our two-stream network. Moreover, the multi-scale feature constructed by convseq2seq [14] ignores the principle of time consistency.

4.5. Ablation Analysis

Evaluation of TST block. In Table 5, we adopt different convolution depths of spatial-temporal modeling blocks and compare their parameter quantity to confirm our proposed depth’s effectiveness. We use residual connections every 5 convolution layers for feature retaining, starting from the first 3*3 convolution layer, noted as a residual block. In the structure of TST, we examine the effect of every residual

Table 5: Ablation experiments on the different depths of the trajectory spatial-temporal (TST) module. TST-11(Ours) shows the best capability of human motion modeling.

Milliseconds	80	160	320	400
TST-6	10.2	23.4	49.1	59.6
TST-16	9.9	22.8	48.8	59.3
TST-21	9.9	22.8	48.7	59.2
TST-11(ours)	9.8	22.6	48.1	58.4

block. Because of the 3*3 convolution layer at the end of each TST block, the experiment is set to ablation experiment of 6, 11, 16 and 21 convolutional layers. As we can see in Table 5, when using 6 convolution layers (TST-6), the result shows a large margin between the 11 convolution layers (TST-11). In this case, we believe the network doesn’t fully capture the global information. However, as the net goes deeper as 16 convolution layers (TST-16) and 21 convolution layers (TST-21), the error doesn’t reduce as expected. Therefore, higher parameter quantity and deeper modeling network don’t always bring an effect on modeling human dynamics.

As mentioned in Methodology, our TST block’s depth is designed by the 22 main joints in H3.6M. The number of joints in 3DPW is 24, which is different from H3.6M. But the joints share nearly the same coordinates at the end of each trunk. Our network can still learn the dynamics of the whole human body. Meanwhile, the number of key joints in CMU-Mocap is 25. When using TST-11, the receptive field of our network doesn’t cover all the 25 joints. However, the joints of CMU-Mocap also share nearly the same coordinates at the end of each trunk. Generally speaking, 11 convolution layers can maximize network performance on different datasets.

Evaluation of TF module. We use direct concatenation instead of temporal fusion for fusing two-stream features from position space and velocity space. When implementing velocity vector without temporal fusion(denoted as TF:× in Table 6), the model shows a higher error on most prediction results. Comparing with TrajectoryNet, which models position space dynamics only, there’s still a gap on average prediction on all time steps. In this case, we believe the features from velocity space disturb the position space modeling when using concatenation. Because of the lack of temporal consistency of the two-stream features through direct concatenation, the network doesn’t model the position and velocity information of the human body at the same time-step when forward modeling, but treats it as a single posture evolution along the time sequence leading to spatial-temporal in conformity. Therefore, abandoning temporal fusion only makes the velocity vector interfere with the prediction result of position space.

Table 6: Influence of temporal fusion on H3.6M. Note that, it leads to poorer performance when using concatenation for two-stream fusion instead of temporal fusion.

	Walking				Eating				Smoking				Discussion			
Milliseconds	80	160	320	400	80	160	320	400	80	160	320	400	80	160	320	400
TrajectoryNet[21]	8.2	14.9	30.0	35.4	8.5	18.4	37.0	44.8	6.3	12.8	23.7	27.8	7.5	20.0	41.3	47.8
TF: \times	8.1	15.1	29.7	34.9	8.3	17.7	36.5	44.3	6.5	12.6	22.3	26.5	7.6	19.6	39.4	46.5
TF: \surd	7.8	14.8	28.9	34.4	7.7	17.1	34.5	42.6	6.2	12.7	23.3	27.6	7.1	18.7	39.7	46.8
	Directions				Greeting				Phoning				Posing			
Milliseconds	80	160	320	400	80	160	320	400	80	160	320	400	80	160	320	400
TrajectoryNet[21]	9.7	22.3	50.2	61.7	12.6	28.1	67.3	80.1	10.7	18.8	37.0	43.1	6.9	21.3	62.9	78.8
TF: \times	10.3	24.8	56.0	68.1	13.3	30.1	74.0	89.7	11.0	19.5	38.5	44.1	7.3	21.7	62.2	78.0
TF: \surd	9.4	22.9	53.7	65.1	12.7	28.0	64.6	79.1	10.0	18.6	37.9	44.4	6.8	21.6	63.5	79.9
	Purchases				Sitting				SittingDown				Taking Photo			
Milliseconds	80	160	320	400	80	160	320	400	80	160	320	400	80	160	320	400
TrajectoryNet[21]	17.1	36.1	64.3	75.1	9.0	22.0	49.4	62.6	10.7	28.8	55.1	62.9	5.4	13.4	36.2	47.0
TF: \times	17.8	37.2	72.4	86.7	9.4	22.4	49.4	63.6	10.7	28.6	56.0	67.1	5.5	13.9	36.4	47.6
TF: \surd	17.4	36.1	59.4	67.4	8.6	21.5	50.3	63.6	11.0	27.4	51.2	60.7	5.5	13.5	37.1	49.0
	Waiting				Walking Dog				Walking Together				Average			
Milliseconds	80	160	320	400	80	160	320	400	80	160	320	400	80	160	320	400
TrajectoryNet[21]	8.2	21.0	53.4	68.9	23.6	52.0	98.1	116.9	8.5	18.5	33.9	43.4	10.2	23.2	49.3	59.7
TF: \times	8.5	20.6	53.5	67.7	22.2	48.3	99.0	117.6	7.8	17.6	34.5	42.2	10.3	23.3	50.6	61.6
TF: \surd	8.3	20.5	51.4	66.3	20.8	47.7	94.1	108.7	7.7	18.1	31.9	41.1	9.8	22.6	48.1	58.4

5. Conclusion

In this paper, we propose a TF module based two-stream architecture that models position stream and velocity stream features for human motion prediction. To ensure the integration of spatial-temporal co-occurrence, the TF module fuses the features from two streams in chronological order, which maintains temporal consistency and shows an effect on feature fusion comparing with related works. Meanwhile, the introduced high-dimensional information, which is the velocity vector, shows its advantage on both short-term modeling and long-term modeling for movement action predicting. Our future work will focus on the TF module generalization that can be adapted to the fusion phase of other two-stream or multi-stream modeling networks.

References

- [1] J. Hamill and K. M. Knutzen. Biomechanical basis of human movement. Lippincott Williams & Wilkins, 2006.
- [2] X. Jin, H. Xiao, X. Shen, J. Yang, Z. Lin, Y. Chen, Z. Jie, J. Feng and S. Yan. Predicting scene parsing and motion dynamics in the future. Advances in Neural Information Processing Systems (NIPS), 2017, pp. 6915–6924.
- [3] Y. T. Xu, Y. Li, and D. Meger. Human Motion Prediction Via Pattern Completion in Latent Representation Space. Conference on Computer & Robot Vision, 2019, DOI: 10.1109/CRV.2019.00016.
- [4] S. Song, C. Lan, J. Xing, W. Zeng and J. Liu. 2017. An end-to-end spatio-temporal attention model for human action recognition from skeleton data. AAAI Conference on Artificial Intelligence, 2017, DOI: arXiv:1611.06067.
- [5] S. Song, C. Lan, J. Xing, W. Zeng and J. Liu. Spatio-Temporal Attention-Based LSTM Networks for 3D Action Recognition and Detection. IEEE Transactions on Image Processing, 2018, vol. 27, no. 7, pp. 3459-3471.
- [6] M. Brand and A. Hertzmann. Style machines. In: Proceedings of the 27th annual conference on Computer graphics and interactive techniques. ACM Press/Addison-Wesley Publishing Co, 2000, pp. 183–192.
- [7] V. Pavlovic, J. M. Rehg and J. MacCormick. Learning switching linear models of human motion. In Advances in neural information processing systems, 2001, pp. 981–987.
- [8] G. E. Hinton, R. R. Salakhutdinov. Reducing the dimensionality of data with neural networks. Science, 2006, 313(5786):7-504.
- [9] G. E. Hinton G E, P. Dayan, B. J. Frey and R. M. Neal. The "wake-sleep" algorithm for unsupervised neural networks. Science, 1995, 268(5214):61-1158.
- [10] D. Sarma, V. Kavyasree and M. K. Bhuyan. Two-stream Fusion Model for Dynamic Hand Gesture Recognition using 3D-CNN and 2D-CNN Optical Flow guided Motion Template. arXiv preprint, 2020, DOI: arXiv: 2007.08847.
- [11] H. Wang and J. Feng. 2019. VRED: A Position-Velocity Recurrent Encoder-Decoder for Human Motion Prediction. arXiv preprint, 2019, DOI: arXiv: abs/1906.06514.
- [12] W. Mao, M. Liu, M. Salzmann and H. Li. Learning trajectory dependencies for human motion prediction. IEEE International Conference on Computer Vision (ICCV), 2019, pp. 9489–9497.
- [13] J. Martinez, M. J. Black and J. Romero. On human motion prediction using recurrent neural networks. IEEE Conference on Computer Vision and Pattern Recognition (CVPR), 2017, pp. 4674-4683.
- [14] C. Li, Z. Zhang, W. S. Lee and G. H. Lee. Convolutional sequence to sequence model for human dynamics. IEEE Conference on Computer Vision and Pattern Recognition, 2018, pp. 5226–5234.
- [15] L. Shi, Y. Zhang and H. Lu. Skeleton-Based Action Recognition with Multi-Stream Adaptive Graph Convolutional Networks. IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR), 2019, DOI: 10.1109/CVPR.2019.00532.
- [16] L. Shi, Y. Zhang, J. Cheng and H. Lu. Two-Stream Adaptive Graph Convolutional Networks for Skeleton-Based Action Recognition. IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR), 2018, pp. 12018-12027.
- [17] A. Gopalakrishnan, A. Mali, D. Kifer, C. L. Giles and A. G. Ororbia.

- A neural temporal model for human motion prediction. *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2019, pp. 12108-12117.
- [18] L. Y. Gui, Y. X. Wang, X. Liang and J. M. Moura. Adversarial geometry-aware human motion prediction. *European Conference on Computer Vision (ECCV)*, 2018, pp. 823-842.
- [19] K. Fragkiadaki, S. Levine, P. Felsen and J. Malik. Recurrent network models for human dynamics. *IEEE International Conference on Computer Vision (ICCV)*, 2015, pp. 4346-4354.
- [20] Z. Liu, S. Wu, S. Jin, Q. Liu, S. Lu, R. Zimmermann and L. Cheng. Towards natural and accurate future motion prediction of humans and animals. *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2019, pp. 9996-10004.
- [21] X. Liu, J. Yin, J. Liu, P. Ding, J. Liu and H. Liub. TrajectoryCNN: a new spatio-temporal feature learning network for human motion prediction. *IEEE Transactions on Circuits and Systems for Video Technology*, 2020, DOI: 1910.06583.
- [22] J. Gehring, M. Auli, D. Grangier, D. Yarats and Y. N. Dauphin. Convolutional sequence to sequence learning. *Proceedings of the 34th International Conference on Machine Learning*, volume 70 of *Proceedings of Machine Learning Research*, 2017, pp. 1243–1252.
- [23] K. Simonyan, A. Zisserman. Two-Stream Convolutional Networks for Action Recognition in Videos[J]. *Advances in neural information processing systems*, 2014, DOI: 10.1002/14651858.CD001941.pub3.
- [24] C. Feichtenhofer, A. Pinz and A. Zisserman. Convolutional Two-Stream Network Fusion for Video Action Recognition. *IEEE Conference on Computer Vision Pattern Recognition*, 2016, DOI: 10.1109/CVPR.2016.213.
- [25] L. Wang, Y. Xiong, Z. Wang, et al. Temporal Segment Networks: Towards Good Practices for Deep Action Recognition. *European Conference on Computer Vision (ECCV)*, 2016, DOI: 10.1007/978-3-319-46484-8_2.
- [26] J. Lin, C. Gan and S. Han. TSM: Temporal Shift Module for Efficient Video Understanding. *IEEE/CVF International Conference on Computer Vision (ICCV)*, 2019, pp. 7082-7092.
- [27] L. Y. Gui, Y. X. Wang, D. Ramanan and J. M. Moura. Few-shot human motion prediction via meta-learning. *European Conference on Computer Vision (ECCV)*, 2018, p. 823.
- [28] X. Shu, L. Zhang, G. J. Qi, W. Liu and J. Tang. Spatiotemporal Co-attention Recurrent Neural Networks for Human-Skeleton Motion Prediction. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2021, DOI: 10.1109/TPAMI.2021.3050918.
- [29] K. Sun, B. Xiao, D. Liu and J. Wang. Deep high-resolution representation learning for human pose estimation. *IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2019, pp. 5686-5696.
- [30] C. Ionescu, D. Papava, V. Olaru and C. Sminchisescu. Human3.6m: Large scale datasets and predictive methods for 3d human sensing in natural environments. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 36, no. 7, 2014. pp. 1325–1339.
- [31] T. von Marcard, R. Henschel, M. J. Black, B. Rosenhahn and G. Pons-Moll. Recovering accurate 3d human pose in the wild using imus and a moving camera. *European Conference on Computer Vision (ECCV)*, 2018, pp. 601–617.