

Space-Time-Separable Graph Convolutional Network for Pose Forecasting

Theodoros Sofianos[†], Alessio Sampieri[†], Luca Franco and Fabio Galasso
Sapienza University of Rome, Italy

Abstract

Human pose forecasting is a complex structured-data sequence-modelling task, which has received increasing attention, also due to numerous potential applications. Research has mainly addressed the temporal dimension as time series and the interaction of human body joints with a kinematic tree or by a graph. This has decoupled the two aspects and leveraged progress from the relevant fields, but it has also limited the understanding of the complex structural joint spatio-temporal dynamics of the human pose. Here we propose a novel Space-Time-Separable Graph Convolutional Network (STS-GCN) for pose forecasting. For the first time, STS-GCN models the human pose dynamics only with a graph convolutional network (GCN), including the temporal evolution and the spatial joint interaction within a single-graph framework, which allows the cross-talk of motion and spatial correlations. Concurrently, STS-GCN is the first space-time-separable GCN: the space-time graph connectivity is factored into space and time affinity matrices, which bottlenecks the space-time cross-talk, while enabling full joint-joint and time-time correlations. Both affinity matrices are learnt end-to-end, which results in connections substantially deviating from the standard kinematic tree and the linear-time time series. In experimental evaluation on three complex, recent and large-scale benchmarks, Human3.6M [24], AMASS [34] and 3DPW [48], STS-GCN outperforms the state-of-the-art, surpassing the current best technique [35] by over 32% in average at the most difficult long-term predictions, while only requiring 1.7% of its parameters. We explain the results qualitatively and illustrate the graph interactions by the factored joint-joint and time-time learnt graph connections.

Our source code is available at:

<https://github.com/FraLuca/STSGCN>

1. Introduction

Forecasting future human poses is the task of modelling the complex structured-sequence of joint spatio-temporal

[†] indicates equal contribution

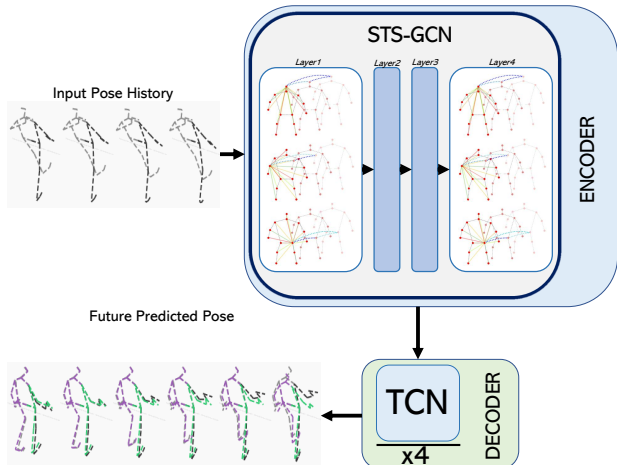


Figure 1: **Overview of the proposed pipeline.** Given a sequence of observed 3D poses, the novel STS-GCN encodes the spatio-temporal body dynamics. The encoded representation serves to predict future poses by means of a Temporal Convolutional Network (TCN). STS-GCN allows the spatial and temporal interaction of joints, cf. green-orange linkage on the Vitruvian man and dashed blue lines connecting joints in time, which are both learnt. But it bottlenecks their cross-talk by a new GCN design with factored space-time adjacency matrices. (Vector image: please zoom in.)

dynamics of the human body. This has received increasing attention due to its manifold applications to autonomous driving [38], healthcare [44], teleoperations [39] and collaborative robots [28, 45], where e.g. anticipating the human motion avoids crashes and helps the robots plan the future.

Research has so far addressed modelling space and time in separate frameworks. Time has generally been modelled with convolutions in the temporal dimension [10], with recurrent neural networks (RNN [35, 36, 49, 11], GRU [53, 1] and LSTM [52]) or with Transformer Networks [9]. Space and the interaction of joints has instead been recently modelled by Graph Convolutional Networks (GCN) [35], mostly connecting body joints along a kine-

matic tree. The separate approach has side-stepped the complexity of a joint model across the spatial and temporal dimensions, which are diverse in nature, and has leveraged progress in the relevant fields. However this has also limited the understanding of the complex human body dynamics.

Here we propose to forecast human motion with a novel Space-Time-Separable Graph Convolutional Network (STS-GCN). STS-GCN encodes both the spatial joint-joint and the temporal time-time correlations with a joint spatio-temporal GCN [27]. The single-graph framework favors the cross-talk of the body joint interactions and their temporal motion patterns. Further to better performance, using the GCN-only model results in considerably less parameters.

To the best of our knowledge, STS-GCN is the first space-time separable GCN. We realize this by factorizing the graph adjacency matrix A^{st} into $A^s A^t$. Our intuition is that bottleneck'ing the cross-talk of the spatial joints and the temporal frames helps to improve the interplay of spatial joints and temporal patterns. This differs substantially from recent work [29, 5] which separate the graph interactions from the channel convolutions, being therefore depthwise separable. Still both separable designs are advantageous for the reduction of model parameters.

Fig. 1 illustrates the encoder-decoder design of our model. Following the body motion encoding by the STS-GCN, the future pose coordinates are forecast with few simple convolutional layers, generally termed Temporal Convolutional Network (TCN) [16, 4, 33], robust and fast to train.

Note from Fig. 1 that the factored $A^s A^t$ graph adjacency matrices are learnt. This results in better performance and it allows us to interpret the joint-joint and the time-time interactions, as we further illustrate in Fig. 3 and in Sec. 4.

In extensive experiments over the modern, challenging and large-scale datasets of Human3.6M [24], AMASS [34] and 3DPW [48], we demonstrate that STS-GCN improves over the state-of-the-art. Notably, STS-GCN outperforms the current best technique [35] by over 32% on all three datasets, in average at the most difficult long-term predictions, while only adopting 1.7% of its parameters.

We summarize our main contributions as follows:

- We propose the first space-time separable graph convolutional network, which is first to factorize the graph adjacency matrix, rather than depthwise [29, 5];
- Our space-time human body representation is the first to exclusively use a GCN and it adopts only 1.7% parameters of the current best competing technique [35];
- We improve on the state-of-the-art by over 32% on Human3.6M [24], AMASS [34] and 3DPW [48], in average at the most challenging long-term predictions;
- The joint-joint and time-time graph edge weights are learnt, which allows to explain their interactions.

2. Related Work

Human pose forecasting is a long-standing problem [10]. We discuss related work by distinguishing the temporal aspects of sequence modelling and the spatial representations. Finally we relate to separable convolutional networks.

Temporal modelling Most recent work in human pose forecasting has leveraged Recurrent Neural Networks (RNN) [15, 25, 37, 11, 35, 36, 49], as well as recurrent variants such as Gated Recurrent Units (GRU) [53, 1] and Long Short-Term Memory Networks [52]. These techniques are flexible, but they have issues with long-term predictions such as inefficient training and poor long-term memory [6, 30, 36, 35]. Research has attempted to tackle this, e.g. by training with generative adversarial networks [18] and by imitation learning [41, 49]. Emerging trends have adopted (self-)attention to model time [35, 9], which also applies to model spatial relations [41, 9].

State-of-the-art performance is also attained with convolutional layers [8, 30, 36, 19, 21, 10] in the temporal dimension, which is known as Temporal Convolutional Networks (TCN) [16, 4, 33]. Here we adopt TCN for future frame prediction due to their performance and robustness, but we encode the space-time body dynamics only with GCN.

Representation of body joints Nearly all literature adopts 3D coordinates or angles. [37] has noted that encoding residuals of coordinates, thus velocity, may be beneficial. [36, 35] has adopted Discrete cosine transform (DCT), thus frequency, which greatly supports for periodic motion. Here we experiment with 3D coordinates and angles, but those representations are compatible with our model.

Representation of human pose Graphs are a natural choice to represent the body. These have mostly been hand-designed, mainly leveraging the natural structure of the kinematic tree [25, 8, 51], and encoded via Graph Convolutional Networks (GCN) [27]. [51] learns the adjacency matrix of the graph, still limiting the connectivity to the kinematic tree. Most recently, research has explored all joints linked together and learnt graph edges [36, 35]. Ours also let the training learn a data-driven graph connectivity and edge weights (see Fig. 3 and Sec. 4 for an illustration).

Separable Convolutions Separable convolutions [40, 13, 22] decouple processing the cross-channel correlations via 1x1 convolutional filters and the spatial correlations via channel-wise spatial convolutions. These are *depthwise*-separable convolutions, based on the hypothesis that the cross-channel and spatial correlations are sufficiently decoupled, so it is preferable not to map them jointly [13].

To the best of our knowledge, only [29] and [5] apply this concept to GCNs, but they design different graph edge weights for different channels, in the spatial [29] or spectral domain [5]. By contrast, our STS-GCN is the first GCN design which separates the graph connectivity itself, by factoring the space-time adjacency matrix. In the spirit of [13],

our hypothesis is that the space-time cross-talk is limited and that decoupling them is more effective and efficient.

3. STS-GCN

The proposed model proceeds by encoding the coordinates of the body joints which are observed in the given input frames and then it leverages the space-time representation to forecast the future joint coordinates. Encoding is modelled by the proposed STS-GCN graph, which considers the interaction of body joints over time, bottleneck'ing the space-time interplay. Decoding future coordinates is modelled with a TCN. In this section, we further provide insights into the STS-GCN model.

3.1. Problem Formalization

We observe the body pose of a person, given by the 3D coordinates or angles of its V joints, for T frames. Then we predict the V body joints for the next K future frames.

We denote the joints by 3D vectors $\mathbf{x}_{v,k}$ representing joint v at time k . The motion history of human poses is denoted by the tensor $\mathcal{X}_{in} = [X_1, X_2, \dots, X_T]$ which we construct out of matrices of 3D coordinates or angles of joints $X_i \in \mathbb{R}^{3 \times V}$ for frames $i = 1 \dots T$. The goal is to predict the future K poses $\mathcal{X}_{out} = [X_{T+1}, X_{T+2}, \dots, X_{T+K}]$.

The motion history tensor is encoded into a graph which models the interaction of all body joints across all observed frames. We define the encoding graph $\mathcal{G} = (\mathcal{V}, \mathcal{E})$, with TV nodes $i \in \mathcal{V}$, which are all body joints across all observed time frames. Edges $(i, j) \in \mathcal{E}$ are represented by a spatio-temporal adjacency matrix $A^{st} \in \mathbb{R}^{VT \times VT}$, relating the interactions of all joints at all times.

3.2. Background on GCN

The spatio-temporal dependencies of joints across times may be conveniently encoded by a GCN, a graph-based neural network model $f(\mathcal{X}_{in}; A, W)$. The input to a graph convolutional layer l is the tensor $\mathcal{H}^{(l)} \in \mathbb{R}^{C^{(l)} \times V \times T}$, which encodes the observed V joints in the T frames. $C^{(l)}$ is the input dimensionality of the hidden representation $\mathcal{H}^{(l)}$. For the first layer, it is $\mathcal{H}^{(1)} = \mathcal{X}_{in}$ and $C^{(1)} = 3$.

A graph convolutional layer l outputs the $\mathcal{H}^{(l+1)} \in \mathbb{R}^{C^{(l+1)} \times V \times T}$, given by the following

$$\mathcal{H}^{(l+1)} = \sigma(A^{st-(l)} \mathcal{H}^{(l)} W^{(l)}) \quad (1)$$

where $A^{st-(l)} \in \mathbb{R}^{VT \times VT}$ is the spatio-temporal adjacency matrix of layer l , $W^{(l)} \in \mathbb{R}^{C^{(l)} \times C^{(l+1)}}$ are the trainable graph convolutional weights of layer l projecting each graph node from $C^{(l)}$ to $C^{(l+1)}$ dimensions, and σ is an activation function such as ReLU, PReLU or tanh.

Two notable graph representations are worth mentioning for their robustness and performance. [51] constrains the graph encoding to the joint-joint relations, thus to a

spatial-only A^s , only along the kinematic tree, and addresses the time-time relations by a convolutional layer of kernel $T \times T \times 1 \times 1$, mapping T frames to T channels. [35], the current state-of-the-art in human pose forecasting, also adopts a spatial-only adjacency matrix A^s , but fully connected. In both cases, the adjacency matrices are trainable.

3.3. Space-Time Separable GCN

The proposed STS-GCN takes motivation from the interaction of the temporal evolution and the spatial joints, as well as from the belief that the interplay of joint-joint and time-time are privileged. Human pose dynamics depend on 3 types of interactions: **i.** joint-joint; **ii.** time-time; and **iii.** joint-time. STS-GCN allows for all 3 types of interactions, but it bottlenecks the joint-time cross-talk.

The interplay of joints over time is modelled by relating the 3 types of relations within a single spatio-temporal encoding GCN. Bottleneck'ing the space-time cross-talk is realized by factoring the space-time adjacency matrix into the product of separate spatial and temporal adjacency matrices $A^{st} = A^s A^t$. A separable space-time graph convolutional layer l is therefore written as follows

$$\mathcal{H}^{(l+1)} = \sigma(A^{s-(l)} A^{t-(l)} \mathcal{H}^{(l)} W^{(l)}) \quad (2)$$

where the same notation as in Eq. (1) applies, apart from the factored $A^{s-(l)} A^{t-(l)}$ of layer l which we explain next.

The adjacency matrix A^s is responsible for the joint-joint interplay. It has dimensionality $A^s \in \mathbb{R}^{V \times V}$, and it models the full joint-joint relations by trainable $V \times V$ matrices for each instant in time (there are T such matrices). Similarly A^t is responsible for the time-time relations. It has dimensionality $A^t \in \mathbb{R}^{T \times T}$ and it defines a full and trainable time-time $T \times T$ relation matrix for each of the V joints.

Note that Eq. 2 represents a single GCN layer, encoding the spatio-temporal interplay of the the body dynamics. The factored space-time matrix bottlenecks the space-time cross-talk, it reduces the model parameters and it yields a considerable increase in the forecasting performance, as we illustrate in Sec. 4. Overall, the graph encoding employs four such GCN layers with residual connections PReLU activation functions, cf. 4 for the implementation details.

Also note that STS-GCN is the sole human pose forecasting graph encoding which exclusively uses GCNs. This contrasts other competing techniques, mostly encoding time with recurrent neural networks [35, 36, 49, 11, 53, 1, 52], or by the use of convolutional layers with kernels across the temporal dimension [51, 10]. This is also a key element to parameter efficiency (see Sec. 4.)

3.4. Discussion on the STS-GCN

Here we first relate STS-GCN to self-attention mechanisms, then we comment on STS-GCN in relation to most recent work on signed and directed GCNs.

msec	Walking				Eating				Smoking				Discussion			
	80	160	320	400	80	160	320	400	80	160	320	400	80	160	320	400
ConvSeq2Seq [30]	17.7	33.5	56.3	63.6	11.0	22.4	40.7	48.4	11.6	22.8	41.3	48.9	17.1	34.5	64.8	77.6
LTD-10-10 [36]	11.1	21.4	37.3	42.9	7.0	14.8	29.8	37.3	7.5	15.5	30.7	37.5	10.8	24.0	52.7	65.8
DCT-RNN-GCN [35]	10.0	19.5	34.2	39.8	6.4	14.0	28.7	36.2	7.0	14.9	29.9	36.4	10.2	23.4	52.1	65.4
Ours	10.7	16.9	29.1	32.9	6.8	11.3	22.6	25.4	7.2	11.6	22.3	25.8	9.8	16.8	33.4	40.2

msec	Directions				Greeting				Phoning				Posing				Purchases				Sitting			
	80	160	320	400	80	160	320	400	80	160	320	400	80	160	320	400	80	160	320	400	80	160	320	400
ConvSeq2Seq [30]	13.5	29.0	57.6	69.7	22.0	45.0	82.0	96.0	13.5	26.6	49.9	59.9	16.9	36.7	75.7	92.9	20.3	41.8	76.5	89.9	13.5	27.0	52.0	63.1
LTD-10-10 [36]	8.0	18.8	43.7	54.9	14.8	31.4	65.3	79.7	9.3	19.1	39.8	49.7	10.9	25.1	59.1	75.9	13.9	30.3	62.2	75.9	9.8	20.5	44.2	55.9
DCT-RNN-GCN [35]	7.4	18.5	44.5	56.5	13.7	30.1	63.8	78.1	8.6	18.3	39.0	49.2	10.2	24.2	58.5	75.8	13.0	29.2	60.4	73.9	9.3	20.1	44.3	56.0
Ours	7.4	13.5	29.2	34.7	12.4	21.8	42.1	49.2	8.2	13.7	26.9	30.9	9.9	18.0	38.2	45.6	11.9	21.3	42.0	48.7	9.1	15.1	29.9	35.0

msec	Sitting Down				Taking Photo				Waiting				Walking Dog				Walking Together				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
ConvSeq2Seq [30]	20.7	40.6	70.4	82.7	12.7	26.0	52.1	63.6	14.6	29.7	58.1	69.7	27.7	53.6	90.7	103.3	15.3	30.4	53.1	61.2	16.6	33.3	61.4	72.7
LTD-10-10 [36]	15.6	31.4	59.1	71.7	8.9	18.9	41.0	51.7	9.2	19.5	43.3	54.4	20.9	40.7	73.6	86.6	9.6	19.4	36.5	44.0	11.2	23.4	47.9	58.9
DCT-RNN-GCN [35]	14.9	30.7	59.1	72.0	8.3	18.4	40.7	51.5	8.7	19.2	43.4	54.9	20.1	40.3	73.3	86.3	8.9	18.4	35.1	41.9	10.4	22.6	47.1	58.3
Ours	14.4	23.7	41.9	47.9	8.2	14.2	29.7	33.6	8.6	14.7	29.6	35.2	17.6	29.4	52.6	59.6	8.6	14.3	26.5	30.5	10.1	17.1	33.1	38.3

Table 1: MPJPE error in mm for short-term prediction of 3D joint positions on Human3.6M. Our model outperforms the state-of-the-art by a large margin. The margin is smaller for very-short-term predictions on periodic actions, e.g. 2-4 frame (80-160 msec) for *Walking* and *Eating*. The margin is larger for the more challenging case of longer-term and aperiodic actions, e.g. up to 40% for *Posing* at 10-frame (400 msec). See Sec. 4.2 for the discussion.

Separable graph convolutions and self-attention Most recent pose forecasting work has leveraged self-attention to encode the relation of frames [35, 9] and/or the relation of joints [9]. Here we relate the proposed STS-GCN to the self-attention mechanisms of [35, 9]. Finally we relate these to Graph Attention Networks (GAT) [47].

Let us first re-write part of the GCN layer of Eq. 1 with the Einstein summation, omitting the indication of layer l , the projection matrix W and the non-linearity σ for better clarity of notation:

$$A^{st}\mathcal{H} = \sum_{vm} A_{wkv}^{st} \mathcal{H}_{vmc} \quad (3)$$

having explicitly indicated with indexes the dimensions of $A^{st} \in \mathbb{R}^{V^T \times V^T}$ and $\mathcal{H} \in \mathbb{R}^{C \times V \times T}$, i.e. indexing spatial joints as $v, w = 1, \dots, V$ and times with $m, k = 1, \dots, T$.

Let us now re-write the corresponding part of the STS-GCN layer of Eq. 2 with the Einstein summation, again omitting the projection matrix W and the non-linearity σ for clarity of notation:

$$A^s(A^t\mathcal{H}) = \sum_v A_{wkv}^s \left(\sum_m A_{kvm}^t \mathcal{H}_{vmc} \right) = \sum_v A_{wkv}^s \mathcal{H}_{kvc}^t \quad (4)$$

where, as above, we have indicated indexes for $A^s \in \mathbb{R}^{V \times V}$ (for each of the T times) and $A^t \in \mathbb{R}^{T \times T}$ (for each of the V joints) as $v, w = 1, \dots, V$ for the spatial joints and as $m, k = 1, \dots, T$ for the times.

Let us now turn to the current best technique for pose forecasting [35]. They adopt a GCN for modelling the spatial interaction of joints at the same time, which coincides with the rightmost term in Eq. 4.

Their temporal modelling is however different from ours, as they adopt an attention formulation $\sigma(QK)V$. Writing it

with the Einstein summation yields:

$$\sum_m \left(\sum_c Q_{kvc}^t K_{vcm}^t \right) V_{vmc} = \sum_m A_{kvm}^{QK-t} V_{vmc} \quad (5)$$

Comparing the right term of Eq. 5 with the separable temporal GCN (the term within parentheses in Eq. 4), we note that the approach of [35], modelling space and time with the different mechanisms of GCN and attention, may also be explained as a separable space-time GCN. The main difference is that A^{QK-t} is a function of the product of inner representation vectors, both stemming from \mathcal{H} . By contrast our temporal adjacency matrix A^t learns the specific pair-wise interaction of relative time shifts. Similar arguments apply when comparing the proposed STS-GCN with the recent GAT [47]. We evaluate the difference *wrt* [35] quantitatively and conduct ablation studies on the adjacency matrices in Sec. 4.

Signed and Directed GCNs Let us now consider that the adjacency matrix A^{st} and its factored terms, A^s and A^t , are trainable parameters. Adjacency matrices were similarly trained by [35], which considers a fully connected matrix, and by [51], which defines specific learnable parameters (denoted M in [51]) to multiply the manually-constructed graph (based on the kinematic tree and the sequential time connections). When encoding the spatio-temporal body dynamics, trainable parameters yield better performance and match the intuition, i.e. they learn the interaction between specific joints and at certain relative temporal offsets.

Trainable parameters result in signed and directed GCNs (see Figs. 3 for an illustration). Both aspects have been surveyed recently [7, 50]. In particular, recent work from [43, 31, 3] maintain that directed graphs encode richer information from their neighborhood, instead of being limited to

distance ranges. Similarly, recent work from [14] demonstrate the superior performance of signed GCNs.

Following the classification of [7, 50], the proposed STS-GCN and the GCNs of [35, 51] are spatial GCNs. This follows from their non-symmetric and possibly ill-posed signed Laplacian matrices, which do not have orthogonal eigendecompositions and are not easily interpretable by spectral-domain constructions [7]. We maintain this makes an interesting direction for future investigation, only partly addressed by very recent work [46].

3.5. Decoding future coordinates

Given the encoded observed body dynamics, the estimation of the 3D coordinates or angles of the body joints in the future is delegated to convolutional layers applying to the temporal dimension. These map the observed frames into the future horizon and refine the estimates via a multi-layered architecture.

Altogether, these layers make a decoder which is generally dubbed Temporal Convolutional Networks (TCN) [16, 4, 33]. While several other sequence modelling options are available, including LSTM [20], GRU [12] and Transformer Networks [17], here we adopt TCNs for their simplicity and robustness, further to satisfactory performance [30].

3.6. Training

The proposed architecture is trained end-to-end supervisedly. Supervision is provided by either of the losses that measure error *wrt* ground truth in terms of Mean Per Joint Position Error (MPJPE) [24, 36] and Mean Angle Error (MAE) [37, 30, 18, 49, 35]. The loss based on MPJPE is:

$$L_{MPJPE} = \frac{1}{V(T+K)} \sum_{k=1}^{T+K} \sum_{v=1}^V \|\hat{\mathbf{x}}_{vk} - \mathbf{x}_{vk}\|_2 \quad (6)$$

where $\hat{\mathbf{x}}_{vk} \in \mathbb{R}^3$ denotes the predicted coordinates of the joint v in the frame k and $\mathbf{x}_{vk} \in \mathbb{R}^3$ is the corresponding ground truth. The loss based on MAE is given by:

$$L_{MAE} = \frac{1}{V(T+K)} \sum_{k=1}^{T+K} \sum_{v=1}^V |\hat{\mathbf{x}}_{vk} - \mathbf{x}_{vk}| \quad (7)$$

where $\hat{\mathbf{x}}_{vk} \in \mathbb{R}^3$ denotes the predicted joint angles in exponential map representation of the joint v in the frame k and $\mathbf{x}_{vk} \in \mathbb{R}^3$ is its ground truth.

4. Experimental evaluation

We experimentally evaluate the proposed model against the state-of-the-art on three recent, large-scale and challenging benchmarks, Human3.6M [24], AMASS [34] and 3DPW [48]. Additionally we conduct ablation studies, evaluate the model qualitatively and illustrate what spatio-temporal graph \mathcal{G} is trained from data.

4.1. Datasets and metrics

Human3.6M [24] The dataset is wide-spread for human pose forecasting and large, consisting of 3.6 million 3D human poses and the corresponding images. It consists of 7 actors performing 15 different actions (e.g. *Walking, Eating, Phoning*). The actors are represented as skeletons of 32 joints. The orientation of joints are represented as exponential maps, from which the 3D coordinates may be computed [42, 15]. For each pose, we consider 22 joints out of the provided 32 for estimating MPJPE and 16 for the MAE. Following the current literature [36, 35, 37], we use the subject 11 (S11) for validation, the subject 5 (S5) for testing, and all the rest of the subjects for training.

AMASS [34] The Archive of Motion Capture as Surface Shapes (AMASS) dataset has been recently proposed, to gather 18 existing mocap datasets. Following [35], we select 13 from those and take 8 for training, 4 for validation and 1 (*BMLrub*) as the test set. Then we use the SMPL [32] parameterization to derive a representation of human pose based on a shape vector, which defines the human skeleton, and its joints rotation angles. We obtain human poses in 3D by applying forward kinematics. Overall, AMASS consists of 40 human-subjects that perform the action of walking. Each human pose is represented by 52 joints, including 22 body joints and 30 hand joints. Here we consider for forecasting the body joints only and discard from those 4 static ones, leading to an 18-joint human pose. As for [24], also these sequences are downsampled to 25 fps.

3DPW [48] The 3D Pose in the Wild dataset [48] consists of video sequences acquired by a moving phone camera. 3DPW includes indoor and outdoor actions. Overall, it contains 51,000 frames captured at 30Hz, divided into 60 video sequences. We use this dataset to test generalization of the models which we train AMASS.

Metrics Following the benchmark protocols, we adopt the MPJPE and MAE error metrics (see Sec. 3.6). The first quantifies the error of the 3D coordinate predictions in mm. The second measures the angle error in degrees. We follow the protocol of [36] and compute MAE with Euler angles. Due to this representation, MAE suffers from an inherent ambiguity, and MPJPE is more effective [9, 2], so mostly adopted here.

Implementation details The graph encoding is given by 4 layers of STS-GCN, which only differ in the number of channels $C^{(l)}$: from 3 (the input 3D coordinates x,y,z or angles), to 64, then 32, 64 and finally 3 (cf. Sec. 3.3), by means of the projection matrices $W^{(l)}$. At each layer we adopt batch normalization [23] and residual connections. Our code is in Pytorch and uses ADAM [26] as optimizer. The learning rate is set to 0.01 and decayed by a factor of 0.1 every 5 epochs after the 20-th. The batch size is 256. On Human3.6M, training for 30 epochs on an NVIDIA RTX 2060 GPU takes 20 minutes.

<i>msec</i>	Walking				Eating				Smoking				Discussion			
	560	720	880	1000	560	720	880	1000	560	720	880	1000	560	720	880	1000
ConvSeq2Seq [30]	72.2	77.2	80.9	82.3	61.3	72.8	81.8	87.1	60.0	69.4	77.2	81.7	98.1	112.9	123.0	129.3
LTD-50-25 [36]	50.7	54.4	57.4	60.3	51.5	62.6	71.3	75.8	50.5	59.3	67.1	72.1	88.9	103.9	113.6	118.5
LTD-10-25 [36]	51.8	56.2	58.9	60.9	50.0	61.1	69.6	74.1	51.3	60.8	68.7	73.6	87.6	103.2	113.1	118.6
LTD-10-10 [36]	53.1	59.9	66.2	70.7	51.1	62.5	72.9	78.6	49.4	59.2	66.9	71.8	88.1	104.4	115.5	121.6
DCT-RNN-GCN [35]	47.4	52.1	55.5	58.1	50.0	61.4	70.6	75.5	47.5	56.6	64.4	69.5	86.6	102.2	113.2	119.8
Ours	40.6	45.0	48.0	51.8	33.9	40.2	46.2	52.4	33.6	39.6	45.4	50.0	53.4	63.6	72.3	78.8

<i>msec</i>	Directions				Greeting				Phoning				Posing				Purchases				Sitting			
	560	720	880	1000	560	720	880	1000	560	720	880	1000	560	720	880	1000	560	720	880	1000	560	720	880	1000
ConvSeq2Seq [30]	86.6	99.8	109.9	115.8	116.9	130.7	142.7	147.3	77.1	92.1	105.5	114.0	122.5	148.8	171.8	187.4	111.3	129.1	143.1	151.5	82.4	98.8	112.4	120.7
LTD-50-25 [36]	74.2	88.1	99.4	105.5	104.8	119.7	132.1	136.8	68.8	83.6	96.8	105.1	110.2	137.8	160.8	174.8	99.2	114.9	127.1	134.9	79.2	96.2	110.3	118.7
LTD-10-25 [36]	76.1	91.0	102.8	108.8	104.3	120.9	134.6	140.2	68.7	84.0	97.2	105.1	109.9	136.8	158.3	171.7	99.4	114.9	127.9	139.9	78.5	95.7	110.0	118.8
LTD-10-10 [36]	72.2	86.7	98.5	105.8	103.7	120.6	134.7	140.9	67.8	83.0	96.4	105.1	107.6	136.1	159.5	175.0	98.3	115.1	130.1	135.9	76.4	93.1	106.9	115.7
DCT-RNN-GCN [35]	73.9	88.2	100.1	106.5	101.9	118.4	132.7	138.8	67.4	82.9	96.5	105.0	107.6	136.8	161.4	178.2	95.6	110.9	125.0	134.2	76.4	93.1	107.0	115.9
Ours	47.6	56.5	64.5	71.0	64.8	76.3	85.5	91.6	41.8	51.1	59.3	66.1	64.3	79.3	94.5	106.4	63.7	74.9	86.2	93.5	47.7	57.0	67.4	75.2

<i>msec</i>	Sitting Down				Taking Photo				Waiting				Walking Dog				Walking Together				Average			
	560	720	880	1000	560	720	880	1000	560	720	880	1000	560	720	880	1000	560	720	880	1000	560	720	880	1000
ConvSeq2Seq [30]	106.5	125.1	139.8	150.3	84.4	102.4	117.7	128.1	87.3	100.3	110.7	117.7	122.4	133.8	151.1	162.4	72.0	77.7	82.9	87.4	90.7	104.7	116.7	124.2
LTD-50-25 [36]	100.2	118.2	133.1	143.8	75.3	93.5	108.4	118.8	77.2	90.6	101.1	108.3	107.8	120.3	136.3	146.4	56.0	60.3	63.1	65.7	79.6	93.6	105.2	112.4
LTD-10-25 [36]	99.5	118.5	133.6	144.1	76.8	95.3	110.3	120.2	75.1	88.7	99.5	106.9	105.8	118.7	132.8	142.2	58.0	63.6	67.0	69.6	79.5	94.0	105.6	112.7
LTD-10-10 [36]	96.2	115.2	130.8	142.2	72.5	90.9	105.9	116.3	73.4	88.2	99.8	107.5	109.7	122.8	139.0	150.1	55.7	61.3	66.4	69.8	78.3	93.3	106.0	114.0
DCT-RNN-GCN [35]	97.0	116.1	132.1	143.6	72.1	90.1	105.5	115.9	74.5	89.0	100.3	108.2	108.2	120.6	135.9	146.9	52.7	57.8	62.0	64.9	77.3	91.8	104.1	112.1
Ours	63.3	73.9	86.2	94.3	47.0	57.4	67.2	76.9	47.3	56.8	66.1	72.0	74.7	85.7	96.2	102.6	38.9	44.0	48.2	51.1	50.8	60.1	68.9	75.6

Table 2: MPJPE error in mm for long-term prediction of 3D joint positions on Human3.6M. Our model outperforms the state-of-the-art by a large margin for each time prediction horizon and each action. Largest improvements *wrt* the current best [35] are obtained for the most challenging cases of longer-term predictions (22-25 frame, 880-1000 msec) of aperiodic actions such as *Sitting* (36%), *Phoning* (43%) and *Posing* (40%). The average improvement over the 14-25 frame (560-1000 msec) predictions is 34%. See Sec. 4.2 for the discussion.

4.2. Comparison to the state-of-the-art

We quantitatively evaluate our proposed model against the state-of-the-art both for short-term (<500 msec) and long-term (>500 msec) predictions.

We include into the comparison: ConvSeq2Seq [30], which adopts convolutional layers, separately encoding long- and short-term history; LTD-X-Y [36], which encodes the sequence frequency with a DCT, prior to a GCN (X and Y stand for the number of observed and predicted frames); BC-WGAIL-div [49], adopting reinforcement learning; and finally DCT-RNN-GCN [35], the current best performer, which extends LTD-X-Y with an RNN and motion-attention.

All algorithms take as input 10 frames (400 msec), with the exception of LTD, for which we also report the case of larger number of input frames. Then algorithms predict future poses for the next 2 to 10 frames (80-400 msec) in the case of short-term, and for 14-25 frames (560-1000 msec) in the long-term case.

Human3.6M: 3D Joint Positions Let us consider Tables 1 and 2 for the tests on short- and long-term prediction respectively. Across all time horizons in both tables, our model outperforms all competing techniques, with the only exception of 3 experiments out of 120 (2-frame predictions for *Walking*, *Eating* and *Directions*), where it is within a marginal error.

Considering the average errors in Table 1, the improvement of our model over the current best [35] ranges from 3% in the case of 2 time frames, up to 34% improvement

for the more challenging case of 10 frames. Note that, at the 10-frame horizon, improvements are less in the case of periodic actions such as *Walking* (17%) but larger for aperiodic actions such as *Posing* (40%). We believe this is because of the DCT encoding of [35].

We illustrate in Table 2 the more arduous long-term prediction horizons. Our predictions at 560 msec (14 frames) are more accurate than those of [35] by 27 mm, while at 1 sec (25 frames) our model reaches an improvement of 37 mm. In average across predictions over 14-25 frames, our model outperforms the current best [35] by 34%.

<i>msec</i>	Average							
	80	160	320	400	560	720	880	1000
LTD-10-25 [36]	0.34	0.57	0.93	1.06	1.27	1.44	1.57	1.66
LTD-10-10 [36]	0.32	0.55	0.91	1.04	1.26	1.44	1.59	1.68
BC-WGAIL-div [49]	0.31	0.57	0.90	1.02	1.23	-	-	1.65
DCT-RNN-GCN [35]	0.31	0.55	0.90	1.04	1.25	1.42	1.56	1.65
Ours	0.24	0.39	0.59	0.66	0.79	0.92	1.00	1.09

Table 3: Average MAE prediction errors over all actions of Human3.6M. Our model improves consistently over the state-of-the-art by a large margin.

Human3.6M: Joint Angles Average angle errors are reported in Table 3. Our model outperforms the current best [35] with larger improvements on the long-term horizon. The performance increase is 23% for 2 frames and it is 34% for 25 future frames.

AMASS Also in the case of AMASS, in Table 4, for short and long-term predictions of 3D coordinates, our model outperforms the state-of-the-art by 32% on the longest time

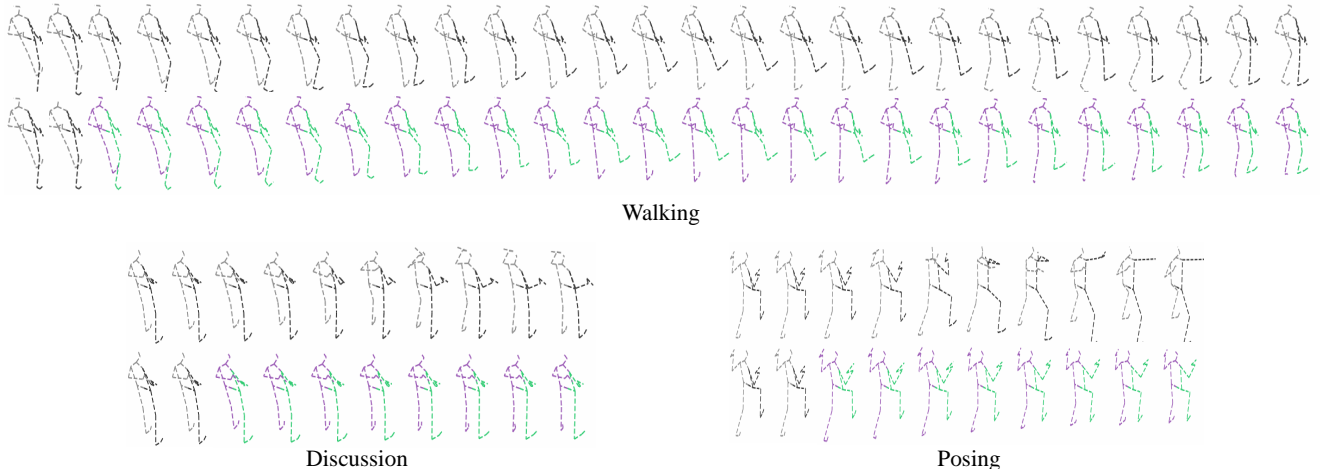


Figure 2: Sample long-term predictions (25 frames, 1 sec) for the actions of *Walking*, *Discussion* and *Posing*. (One every three frames are shown for the second two.) Purple/green limbs are the left/right sides of the body. Gray/black pictorials indicate the observed ground-truth (GT) skeletons. Predictions accurately match the GT. Mistakes may be observed on the left hand of the person in *Discussion* and in the aperiodic motion of *Posing*, an action performed in different ways in the training dataset. Zoom in for details.

msec	AMASS-BMLrub							
	80	160	240	400	560	720	880	1000
convSeq2Seq [30]	20.6	39.6	59.7	67.6	79.0	87.0	91.5	93.5
LTD-10-10 [36]	10.3	19.3	36.6	44.6	61.5	75.9	86.2	91.2
LTD-10-25 [36]	11.0	20.7	37.8	45.3	57.2	65.7	71.3	75.2
DCT-RNN-GCN [35]	11.3	20.7	35.7	42.0	51.7	58.6	63.4	67.2
Ours	10.0	12.5	21.8	24.5	31.9	38.1	42.7	45.5

Table 4: Average MPJPE in mm over the *BMLrub* test sequences of AMASS. Our model outperforms the current best [35] by 32% for 25-frame (1000 msec) predictions.

horizon (25-frame, 1000 msec).

3DPW In Table 5, we test the generalizability of our model by training on AMASS and testing on 3DPW. Results are significantly beyond the state-of-the-art. For 2-frame predictions we reduce the error by 32%, compared to the second best. For any other time horizon above 4 frames, we reduce the error by at least 43%.

Number of parameters Table 6 (*rightmost column*) compares the number of parameters of our model Vs. [35]. Ours uses a fraction of parameters, 57.5k Vs 3.4M, only 1.7%.

Qualitative evaluation We provide sample predictions (*purple/green*) in Fig. 2 on Human3.6M against ground truth sequences (*gray/black*). All predictions are long-term (25 frames) but we only display one every three frames for *Discussion* and *Posing*, to fit the illustrations into a row. Results are in line with the long-term error statistics of Table 2. The forecast *Walking* is accurate, within 5.2 cm-accuracy in average at 25 frames (1 sec) and pictorially matching the

msec	3DPW							
	80	160	240	400	560	720	880	1000
convSeq2Seq [30]	18.8	32.9	52.0	58.8	69.4	77.0	83.6	87.8
LTD-10-10 [36]	12.0	22.0	38.9	46.2	59.1	69.1	76.5	81.1
LTD-10-25 [36]	12.6	23.2	39.7	46.6	57.9	65.8	71.5	75.5
DCT-RNN-GCN [35]	12.6	23.1	39.0	45.4	56.0	63.6	69.7	73.7
Ours	8.6	12.8	21.0	24.5	30.4	35.7	39.6	42.3

Table 5: Average MPJPE in mm, testing the generalizability on 3DPW of models trained on AMASS. Our model scores significantly beyond the state-of-the-art, i.e. it outperforms [35], on 4-25 frames (160-1000 msec) by at least 43%.

ground truth. This shows how our model learns periodic motion well. Predicted future poses are also relatively accurate for *Discussion*, where the average error is 7.9 cm (cf. competing algorithms are nearly 12 cm). In this case, our model predicts well the mostly static pose of the discussing person, but the error is larger on the waving left hand. Finally our model is producing larger errors on *Posing* (10.6 cm in average), as it is a more challenging aperiodic action, which different people perform in different ways.

4.3. Ablation Study

Table 6 illustrates the following ablative variants of our proposed STS-GCN encoding technique:

Distinct graphs \mathcal{G}^s and \mathcal{G}^t This stands for separate GCNs for space and time, with separate adjacency and projection matrices, intertwined by an activation function. The variant underperforms our proposed model, which confirms the

msec	Average								Parameters
	80	160	240	400	560	720	880	1000	
DCT-RNN-GCN [35]	10.4	22.6	47.1	58.3	77.3	91.8	104.1	112.1	3.4M
Distinct $\mathcal{G}^s, \mathcal{G}^t$	28.9	26.4	40.2	48.7	58.7	66.9	75.2	79.9	59.8k
Full \mathcal{G}^{st}	11.9	19.4	34.1	40.8	53.1	65.6	75.1	82.5	222.9k
Separable \mathcal{G}^{s-t} shared	11.3	19.4	34.7	40.5	52.5	62.1	69.2	76.9	36.4k
Separable \mathcal{G}^{s-t} (proposed)	10.1	17.1	33.1	38.3	50.8	60.1	68.9	75.6	57.5k

Table 6: Average MPJPE error in mm on Human3.6M, comparing ablating variants of our model. See 4.3 for the detailed discussion. We also report here the number of parameters of all techniques, as well as of the current best algorithm [35]. Our proposed Separable \mathcal{G}^{s-t} has only 1.7% of the parameters of [35].

importance of spatio-temporal interaction within a single graph \mathcal{G}^{st} . Interestingly the errors are much larger for the short-range (nearly 3x larger) than for the long-range (+6% errors). We believe longer-term correlations may aid the variant.

Full (non-separable) graph \mathcal{G}^{st} The variant adopts a full space-time adjacency matrix A^{st} . We observe a similar trend as for distinct graphs, i.e. worse performance with larger error increase (+18%) for short-term predictions but better for long-term ones (+9%). Notably the full graph model requires nearly 4x more parameters than our proposed one, cf. rightmost column in Table 6.

Separable graph \mathcal{G}^{s-t} shared across layers This only differs as it learns shared adjacency matrices across all layers, rather than layer-specific ones. Errors are comparable in the long-term (+2%) but larger in the short-term (+12%), against saving 37% of the parameters.

Learnt separable space and time adjacency matrices

In Fig. 3, we illustrate two learnt adjacency matrices, upon training on Human3.6M. On the left, we represent a spatial adjacency matrix A^s , i.e. imagine the red dots positioned on the 22 keypoints of a frontally posing Vitruvian man. Learnt parameters are directed (as the learnt matrix A^s is not symmetric) and signed edges (cf. Sec. 3.3), color-coded weights as in the legend. For clarity of illustration, we represent the two strongest connections for each keypoint. Note how most learnt connections follow the kinematic tree, which confirms the importance of the physical linkage. However additional strong connections also emerge, which bridge distant but motion-related joints, such as the two feet, the feet to the head, and the shoulders to the opposite hips, which intuitively interact for future pose prediction.

In Fig. 3 (right), we represent a temporal adjacency matrix A^t , also asymmetric and signed. It is noticeable the information flow from the earlier to the later observed frames. So the bottom-left side of the matrix shows larger absolute values. In particular most information is drawn to the last two frames (bottom two rows), corresponding to the 9th and 10th observed frames. Note also that the range of tempo-

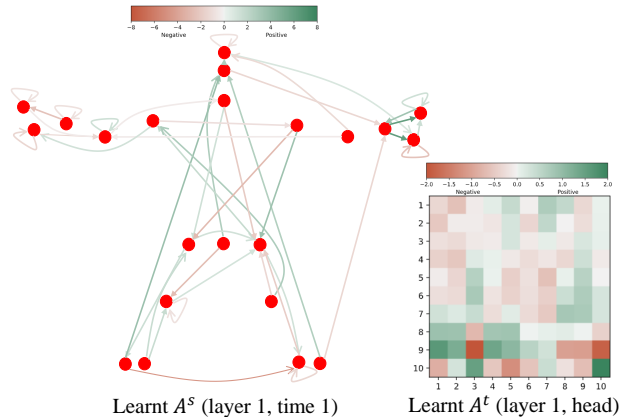


Figure 3: Sample spatial A^s (left) and temporal A^t (right) adjacency matrices, learnt on Human3.6M. (left) Red dots represent the 22 human joints; learnt joint-joint relations mainly follow the kinematic tree, but additionally bring up long-term connections (e.g. foot-foot, head-foot) which support forecasting. (right) The learnt A^t shows information flow from the earlier to the later observed frames, i.e. larger absolute values in the bottom-left part of the matrix. (Zoom in for details.) See Sec. 4.3 for a discussion.

ral relation coefficients $[-2, 2]$ is smaller than the spatial $[-8, 8]$, which privileges spatial information above the temporal when forecasting future poses.

5. Conclusions

We have proposed a novel Space-Time-Separable Graph Convolutional Network (STS-GCN) for pose forecasting. The single-graph framework favors the cross-talk of space and time, while bottleneck’ing the space-time interaction allows to better learn the fully-trainable joint-joint and time-time interactions. The model improves considerably on the state-of-the-art performance and but only requires a fractions of the parameters. These results further support the adoption of GCN and future research on it.

Acknowledgements

The authors wish to acknowledge Panasonic for partially supporting this work and the project of the Italian Ministry of Education, Universities and Research (MIUR) “Dipartimenti di Eccellenza 2018-2022”.

References

- [1] Vida Adeli, Ehsan Adeli, Ian Reid, Juan Carlos Niebles, and Hamid Rezatofighi. Socially and contextually aware human motion and pose forecasting. *IEEE Robotics and Automation Letters (RA-L)*, PP:1–1, 2020. 1, 2, 3
- [2] Ijaz Akhter et al. Pose-conditioned joint angle limits for 3d human pose reconstruction. In *CVPR*, 2015. 5
- [3] James Atwood and Don Towsley. Diffusion-convolutional neural networks. In *Advances in Neural Information Processing Systems (NIPS)*, 2016. 4
- [4] Shaojie Bai, J. Z. Kolter, and V. Koltun. An empirical evaluation of generic convolutional and recurrent networks for sequence modeling. *ArXiv*, abs/1803.01271, 2018. 2, 5
- [5] Muhammet Balcilar, Guillaume Renton, Pierre Héroux, Benoit Gauzère, Sébastien Adam, and Paul Honeine. Spectral-designed depthwise separable graph neural networks. In *The International Conference on Machine Learning (ICML) - Workshop on Graph Representation Learning and Beyond (GRL+ 2020)*, 2020. 2
- [6] Emad Barsoum, John Kender, and Zicheng Liu. HP-GAN: probabilistic 3d human motion prediction via GAN. *Computing Research Repository (CoRR)*, 2017. 2
- [7] M. M. Bronstein, J. Bruna, Y. LeCun, A. Szlam, and P. Vandergheynst. Geometric deep learning: Going beyond euclidean data. *IEEE Signal Processing Magazine*, 34(4):18–42, 2017. 4, 5
- [8] Judith Bütetage, Michael J. Black, D. Kragic, and H. Kjellström. Deep representation learning for human motion prediction and classification. *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 1591–1599, 2017. 2
- [9] Yujun Cai, Lin Huang, Yiwei Wang, Tat-Jen Cham, Jianfei Cai, Junsong Yuan, Jun Liu, Xu Yang, Yiheng Zhu, Xiaohui Shen, Ding Liu, Jing Liu, and Nadia Magnenat Thalmann. Learning progressive joint propagation for human motion prediction. In *The European Conference on Computer Vision (ECCV)*, 2020. 1, 2, 4, 5
- [10] Zhe Cao, Hang Gao, Karttikeya Mangalam, Qizhi Cai, Minh Vo, and Jitendra Malik. Long-term human motion prediction with scene context. In *The European Conference on Computer Vision (ECCV)*, 2020. 1, 2, 3
- [11] H. Chiu, E. Adeli, B. Wang, D. Huang, and J. C. Niebles. Action-agnostic human pose forecasting. In *The IEEE Winter Conference on Applications of Computer Vision (WACV)*, 2019. 1, 2, 3
- [12] Kyunghyun Cho, Bart van Merriënboer, Caglar Gulcehre, Dzmitry Bahdanau, Fethi Bougares, Holger Schwenk, and Yoshua Bengio. Learning phrase representations using RNN encoder–decoder for statistical machine translation. In *The Conference on Empirical Methods in Natural Language Processing (EMNLP)*, 2014. 5
- [13] F. Chollet. Xception: Deep learning with depthwise separable convolutions. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2017. 2
- [14] Tyler Derr, Yao Ma, and Jiliang Tang. Signed graph convolutional network. *arXiv*, abs/1808.06354, 2018. 5
- [15] Katerina Fragkiadaki, Sergey Levine, Panna Felsen, and Jitendra Malik. Recurrent network models for human dynamics. In *The IEEE International Conference on Computer Vision (ICCV)*, 2015. 2, 5
- [16] Jonas Gehring, Michael Auli, David Grangier, Denis Yarats, and Yann N. Dauphin. Convolutional sequence to sequence learning. In *The International Conference on Machine Learning (ICML)*, 2017. 2, 5
- [17] Francesco Giuliari, Irtiza Hasan, Marco Cristani, and Fabio Galasso. Transformer networks for trajectory forecasting. In *The International Conference on Pattern Recognition (ICPR)*, 2020. 5
- [18] Liang-Yan Gui, Yu-Xiong Wang, Xiaodan Liang, and José M. F. Moura. Adversarial geometry-aware human motion prediction. In *The European Conference on Computer Vision (ECCV)*, 2018. 2, 5
- [19] A. Hernandez, J. Gall, and F. Moreno. Human motion prediction via spatio-temporal inpainting. In *The IEEE/CVF International Conference on Computer Vision (ICCV)*, 2019. 2
- [20] Sepp Hochreiter and Jürgen Schmidhuber. Long short-term memory. *Neural Computation*, 9(8):1735–1780, 1997. 5
- [21] Daniel Holden, Jun Saito, Taku Komura, and Thomas Joyce. Learning motion manifolds with convolutional autoencoders. In *SIGGRAPH Asia 2015 Technical Briefs*, 2015. 2
- [22] Andrew G. Howard, Menglong Zhu, Bo Chen, Dmitry Kalenichenko, Weijun Wang, Tobias Weyand, Marco Andreetto, and Hartwig Adam. Mobilenets: Efficient convolutional neural networks for mobile vision applications. *Computing Research Repository (CoRR)*, abs/1704.04861, 2017. 2
- [23] Sergey Ioffe and Christian Szegedy. Batch normalization: Accelerating deep network training by reducing internal covariate shift. In *The International Conference on Machine Learning (ICML)*, 2015. 5
- [24] 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 (TPAMI)*, 2014. 1, 2, 5
- [25] A. Jain, A. R. Zamir, S. Savarese, and A. Saxena. Structural-rnn: Deep learning on spatio-temporal graphs. In *The IEEE Computer Vision and Pattern Recognition (CVPR)*, 2016. 2
- [26] Diederik P. Kingma and Jimmy Ba. Adam: A method for stochastic optimization. In *The International Conference on Learning Representations (ICLR)*, 2015. 5
- [27] Thomas N. Kipf and Max Welling. Semi-supervised classification with graph convolutional networks. In *The International Conference on Learning Representations (ICLR)*, 2017. 2
- [28] H. S. Koppula and A. Saxena. Anticipating human activities for reactive robotic response. In *The IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, pages 2071–2071, 2013. 1
- [29] Guokun Lai, Hanxiao Liu, and Yiming Yang. Learning depthwise separable graph convolution from data manifold.

- In *The International Conference on Learning Representations (ICLR) rej.*, 2018. 2
- [30] Chen Li, Zhen Zhang, Wee Sun Lee, and Gim Hee Lee. Convolutional sequence to sequence model for human dynamics. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2018. 2, 4, 5, 6, 7
- [31] Yaguang Li, Rose Yu, Cyrus Shahabi, and Yan Liu. Diffusion convolutional recurrent neural network: Data-driven traffic forecasting. In *The International Conference on Learning Representations (ICLR)*, 2018. 4
- [32] Matthew Loper, Naureen Mahmood, Javier Romero, Gerard Pons-Moll, and Michael J. Black. Smpl: A skinned multi-person linear model. *ACM Trans. Graph.*, 34(6), 2015. 5
- [33] W. Luo, B. Yang, and R. Urtasun. Fast and furious: Real time end-to-end 3d detection, tracking and motion forecasting with a single convolutional net. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2018. 2, 5
- [34] Naureen Mahmood, Nima Ghorbani, Nikolaus F. Troje, Gerard Pons-Moll, and Michael J. Black. Amass: Archive of motion capture as surface shapes. In *The IEEE International Conference on Computer Vision (ICCV)*, 2019. 1, 2, 5
- [35] Wei Mao, Miaomiao Liu, and Mathieu Salzmann. History repeats itself: Human motion prediction via motion attention. In *The European Conference on Computer Vision (ECCV)*, 2020. 1, 2, 3, 4, 5, 6, 7, 8
- [36] Wei Mao, Miaomiao Liu, Mathieu Salzmann, and Hongdong Li. Learning trajectory dependencies for human motion prediction. In *The IEEE International Conference on Computer Vision (ICCV)*, 2019. 1, 2, 3, 4, 5, 6, 7
- [37] Julieta Martinez, Michael J. Black, and Javier Romero. On human motion prediction using recurrent neural networks. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2017. 2, 5
- [38] B. Paden, M. Čáp, S. Z. Yong, D. Yershov, and E. Frazzoli. A survey of motion planning and control techniques for self-driving urban vehicles. *IEEE Transactions on Intelligent Vehicles (T-IV)*, 1(1):33–55, 2016. 1
- [39] Matteo Rubagotti, Tasbolat Taunyazov, Bukeikhan Omarali, and Almas Shintemirov. Semi-autonomous robot teleoperation with obstacle avoidance via model predictive control. *IEEE Robotics and Automation Letters (RA-L)*, PP:1–1, 2019. 1
- [40] C. Szegedy, Wei Liu, Yangqing Jia, P. Sermanet, S. Reed, D. Anguelov, D. Erhan, V. Vanhoucke, and A. Rabinovich. Going deeper with convolutions. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 1–9, 2015. 2
- [41] Yongyi Tang, Lin Ma, Wei Liu, and Wei-Shi Zheng. Long-term human motion prediction by modeling motion context and enhancing motion dynamic. In *The International Joint Conference on Artificial Intelligence (IJCAI)*, 2018. 2
- [42] Graham W Taylor, Geoffrey E Hinton, and Sam Roweis. Modeling human motion using binary latent variables. In *Advances in Neural Information Processing Systems*, 2007. 5
- [43] Zekun Tong, Yuxuan Liang, Changsheng Sun, David S. Rosenblum, and Andrew Lim. Directed graph convolutional network. *arXiv*, abs/2004.13970, 2020. 4
- [44] Nikolaus Troje. Decomposing biological motion: A framework for analysis and synthesis of human gait patterns. *Journal of vision*, 2:371–87, 2002. 1
- [45] V. V. Unhelkar, P. A. Lasota, Q. Tyroller, R. Buhai, L. Marceau, B. Deml, and J. A. Shah. Human-aware robotic assistant for collaborative assembly: Integrating human motion prediction with planning in time. *IEEE Robotics and Automation Letters (RA-L)*, 3(3):2394–2401, 2018. 1
- [46] J. J. P. Veerman and Robert Lyons. A primer on laplacian dynamics in directed graphs. *arXiv*, abs/2002.02605, 2020. 5
- [47] Petar Veličković, Guillem Cucurull, Arantxa Casanova, Adriana Romero, Pietro Lio, and Y. Bengio. Graph attention networks. In *The International Conference on Learning Representations (ICLR)*, 2018. 4
- [48] Timo von Marcard, Roberto Henschel, Michael Black, Bodo Rosenhahn, and Gerard Pons-Moll. Recovering accurate 3d human pose in the wild using imus and a moving camera. In *The European Conference on Computer Vision (ECCV)*, 2018. 1, 2, 5
- [49] B. Wang, E. Adeli, H. Chiu, D. Huang, and J. C. Niebles. Imitation learning for human pose prediction. In *The IEEE International Conference on Computer Vision (ICCV)*, 2019. 1, 2, 3, 5, 6
- [50] Z. Wu, S. Pan, F. Chen, G. Long, C. Zhang, and P. S. Yu. A comprehensive survey on graph neural networks. *IEEE Transactions on Neural Networks and Learning Systems*, 32(1):4–24, 2021. 4, 5
- [51] Sijie Yan, Yuanjun Xiong, and Dahua Lin. Spatial temporal graph convolutional networks for skeleton-based action recognition. In *The AAAI Conference on Artificial Intelligence*, 2018. 2, 3, 4, 5
- [52] Ye Yuan and Kris Kitani. Ego-pose estimation and forecasting as real-time pd control. In *The IEEE International Conference on Computer Vision (ICCV)*, 2019. 1, 2, 3
- [53] Ye Yuan and Kris Kitani. Dlow: Diversifying latent flows for diverse human motion prediction. In *The European Conference on Computer Vision (ECCV)*, 2020. 1, 2, 3