

Graph-based Normalizing Flow for Human Motion Generation and Reconstruction

Wenjie Yin¹, Hang Yin¹, Danica Kragic¹, Mårten Björkman¹

Abstract—Data-driven approaches for modeling human skeletal motion have found various applications in interactive media and social robotics. Challenges remain in these fields for generating high-fidelity samples and robustly reconstructing motion from imperfect input data, due to e.g. missed marker detection. In this paper, we propose a probabilistic generative model to synthesize and reconstruct long horizon motion sequences conditioned on past information and control signals, such as the path along which an individual is moving. Our method adapts the existing work MoGlow by introducing a new graph-based model. The model leverages the spatial-temporal graph convolutional network (ST-GCN) to effectively capture the spatial structure and temporal correlation of skeletal motion data at multiple scales. We evaluate the models on a mixture of motion capture datasets of human locomotion with foot-step and bone-length analysis. The results demonstrate the advantages of our model in reconstructing missing markers and achieving comparable results on generating realistic future poses. When the inputs are imperfect, our model shows improvements on robustness of generation.

I. INTRODUCTION

Modeling human motion is essential and challenging in human-robot interaction. Computational models that capture rich motion patterns can facilitate animating synthetic characters [1] and understanding human behaviors for greater autonomy in social robotics scenarios [2]. Recent work [3], [4], [5] have proposed deterministic data-driven motion synthesis methods. These methods are often limited to generation of stereotypical samples and fail to capture the natural variability of human motion. Probabilistic generative models, however, permit modelling of the full space of possible poses without collapsing to an average pose [6], [7], [8]. Normalizing flow based approaches have so far received less attention compared to the alternatives, and have rarely been explored for human motion synthesis [1].

One central challenge that generative models are facing is to achieve robust synthesis under imperfect conditions. Human skeletal motion is commonly collected using motion capture (MoCap) systems such as Vicon² or OptiTrack³. When moved to less controlled environments, these systems inevitably suffer from missed markers due to occlusion or other detection failures [9]. Unfortunately, most previous works do not yield satisfactory performance in generating stable and consistent motion patterns under such conditions.

We propose a graph-based probabilistic generative model to address this limitation in the context of human motion generation and reconstruction, as illustrated in Fig. 1. The model builds upon MoGlow [1], an autoregressive normalizing flow

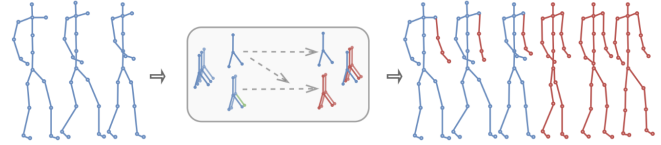


Fig. 1: A probabilistic framework for human motion generation and reconstruction: Generating future human motion given partially observed past pose sequences using graph-based normalizing flow and reconstructing the missing markers in past sequences.

model. The idea is to exploit the invariant spatial correlation of human skeletons and to enforce this inductive bias with a graph structure. A spatial configuration partition strategy orders the nodes in a spatial graph neighbor set, allowing for convolutions on graphs. By using graph convolutional networks, the model encodes correlations between markers and hence achieves a robust synthesis. Our framework is designed as a two-step pipeline. The first step is to generate future poses given incomplete past sequences and control information. The second step is to reconstruct the missing markers by reversing the generated sequences. We evaluate our framework on a mixture of datasets of human locomotion. The evaluation shows that, even with imperfect input data, our proposed graph motion glow model attains a reconstructed motion quality close to that of recorded motion capture data, outperforming a state-of-the-art baseline.

In summary, our work contributes

- an integration of graph neural network and normalizing flows to improve the encoding of human skeletal motion, and
- extending MoGlow to generate human motion patterns given incomplete past poses and managing to reconstruct missing markers robustly.

The paper proceeds as follows: Section II gives an introduction to previous works on the topics of human motion generation and reconstruction, including graph-based human body representation. Section III formulates the problem and describes the full design of our proposed framework for motion generation and reconstruction. Section IV presents the experimental setup and discusses results, compares the results to the baseline. Finally, section V outlines a plan for future work and concludes the paper.

¹Robotics, Perception and Learning lab, KTH Royal Institute of Technology, Sweden. {yinw, hyin, dani, celle}@kth.se.

²<http://www.vicon.com/>

³<https://optitrack.com/>

II. RELATED WORK

In this section, we provide an overview of deep learning based motion generation (Section II-A) and pose reconstruction methods from partially observable data (Section II-B), and then we describe prior works that use graph-based human motion representation (Section II-C).

A. Deep Learning for Motion Prediction and Generation

Deep learning approaches excel at processing massive amounts of data with no requirements on feature engineering. Such approaches have been widely adopted for human motion generation following earlier successes in other domains. Fragkiadaki et al. [3] present an encoder-recurrent-decoder (ERD) to jointly learn a skeleton embedding with sequential information for human pose prediction. Butepage et al. [4] propose a temporal encoding-decoding network to encode previous frames and decode future sequences. Martinez et al. [5] develop a sequence-to-sequence architecture with residual connections to predict joint velocities, while Li et al. [10] propose the auto-conditioned LSTM to synthesize complex and long-term motion patterns.

The approaches reviewed above predict deterministic samples for a given input. To synthesize diverse motion patterns, one way is to use probabilistic generative approaches. Generative adversarial networks (GANs) and their variants have been applied in human motion synthesis and control [11], [6] and speech-driven motion generation [12]. GANs are deemed more effective for encoding modes of the data but can be difficult to train and evaluate. As an alternative, generative autoregressive networks have the advantage of simplicity and have been applied for audio-driven motion generation [7], [13]. Another large branch of methods are based on variational autoencoders (VAEs), which optimize a lower bound on the data log-likelihood. Habibie et al. [8] proposed a VAE with a recurrent structure for controlled motion synthesis. In [14], a hierarchical recurrent model is proposed with each motion sub-sequence mapped to a stochastic latent code through a VAE. VAE-based methods have also found applications in cross-modal synthesis, generating motion from speech [15].

Compared to GANs and VAEs, flow-based generative models have been less explored but are gaining popularity in recent human motion generation works. Henter et al. [1] proposed a normalizing flow-based model called MoGlow that extend Glow [16] to skeleton data. This model is adapted in [17] for an application of speech-driven gesture synthesis. Our approach is an extension to MoGlow. As a result, it features the benefits of flow-based models, which allow for tractable likelihood evaluation and efficient parameterization of encoder and decoder. Flow-based models such as MoGlow typically process data in the Euclidean space and lack the structure for skeleton data with explicit spatial correlations. This may be particularly problematic when strong generalization is expected to cope with untrained uncertainties, e.g. input poses with missing markers. Our method remedies this by integrating graph convolutional networks.

B. Pose Reconstruction from Partially Observable Data

Reconstructing from partial observations has been a long-lasting problem and investigated in human pose estimation. Traditional reconstruction approaches are largely based on matrix factorization. Peng et al. [18] use adaptive non-negative matrix factorization with hierarchical blocks for motion recovery. Wang et al. [19] decompose the entire pose to partial models to exploit the abundant local body posture. Dictionary learning is designed and applied in parallel for each part. Gloersen and Federolf [20] exploit marker inter-correlations from weighted principal components analysis (PCA) for reconstruction. Low-rank matrix completion can be applied to motion recovery [21], [22]. Cui et al. [23] proposed a nonlocal low-rank regularization model (NLR) utilizing kinematic information and weighted Schatten p -norm (WSN) to recover the missing markers. Nevertheless, these approaches make a strong assumption that perfect pose frames should be present in the sequence at least once. In our work, this assumption is not necessary, since in each frame some markers are allowed to be missing.

Deep neural networks, especially recurrent ones, have also been explored to encode temporal correlation for reconstructing missing markers. Mall et al. [24] propose a deep bidirectional LSTM for denoising and synthesizing missing frames. Kucherenko et al. [9] use similar but simpler LSTM-based and time-window-based models. Most methods are variants of LSTM but suffer from the difficulty of capturing long-term dependencies. A deep bi-directional attention network (BAN) [25] is proposed and embedded in the bidirectional LSTM to alleviate this issue. Lohit and Anirudh [26] consider the problem of reconstructing completely unobserved markers of motion sequences. The reconstruction is solved by projecting the observed action onto the action manifold via latent space optimization. Similar to previous approaches, we also use an LSTM to account for temporal causalities. Furthermore, our method encodes spatial and temporal relations between markers and input frames with graph convolutional networks.

C. Graph-based Human Body Representation

Graph neural networks (GNNs) have received increasing attention and have been successfully applied to represent human skeleton data. Most works focus on discriminative tasks such as skeleton-based action recognition [27], [28], [29] and group behavior recognition [30]. Si et al. [31] employ a graph convolutional LSTM network with an attention mechanism to enhance information extraction. Inspired by deformable part-based models (DPMs), [32] divides the skeleton graph into subgraphs and proposes a part-based graph convolutional network (PB-GCN) for action recognition. In [33], motif-based graph convolution is proposed to learn hierarchical spatial structures for action recognition.

For human motion prediction and generation, Li et al. [34] propose a multi-scale graph computational unit (MGCU) to extract deep features. Jain et al. [35] construct a structural graph in which the nodes and edges consist of LSTMs to model the body dynamics. Yan et al. [11] design a graph-based framework called convolutional sequence generation

network (CSGN) to directly generate the entire sequence instead of sequentially. In addition, graph convolution networks have been applied to pose regression [36], trajectory generation [37], and human video prediction [38]. In [39], the framework of normalizing flows is extended with a graph auto-encoder to generate non-human graph structure samples. Inspired by the ST-GCN [27], graph normalizing flows (GNFs) [39] and MoGlow [1], we use spatial graph convolutional network (S-GCN) to encode the spatial relationship of the human skeleton in normalizing flows and leverage ST-GCN to extract features from past sequences.

III. METHOD

This section formulates our target problem and establishes notations used throughout the paper. Preliminaries about normalizing flow, MoGlow and spatial-temporal graph convolutional networks are also given. On the basis of these, we present the contributed framework.

A. Problem Formulation

Suppose that a 3D skeleton-based pose at time t is denoted as $X^{(t)} \in \mathbb{R}^{M \times C}$, with M joint markers and $C = 3$ feature dimensions. The past human poses till time step t_0 are $\mathbb{X}_{(t_0-T_h):(t_0-1)} = [X^{(t_0-T_h)}, \dots, X^{(t_0-1)}] \in \mathbb{R}^{M \times C \times T_h}$. The goal is to estimate a parameterized probabilistic model p_θ from a set of human pose trajectory data, where the optimal parameters are given by:

$$\theta^* = \operatorname{argmax}_\theta p_\theta(\mathbb{X}_{(t_0):(T)} | \mathbb{X}_{(t_0-T_h):(t_0-1)}, C_{(t_0-T_h):(T)}) \quad (1)$$

from which one can sample to predict $T + 1$ future poses $\mathbb{X}_{(t_0):(T)}$, given the past poses $\mathbb{X}_{(t_0-T_h):(t_0-1)}$ and a control input $C_{(t_0-T_h):(T)}$ for the full sequence. Sometimes, the past poses are only partially observed, e.g. with frames missing some marker positions. In such cases, the task becomes that of reconstructing a full trajectory $\mathbb{X}_{(t_0):(T)}$ from an imperfect input, an input that we denote $\hat{\mathbb{X}}_{(t_0-T_h):(t_0-1)}$.

B. Normalizing Flows and MoGlow

Normalizing flows [40], [41], [16] are a class of generative models that allow efficient sampling and inference. The idea is to find an invertible transformation $z = f(X)$ with inverse $X = f^{-1}(f(X))$ to map the data X into a latent space where the distribution is tractable, such as a multivariate Gaussian $p_\theta(z)$. For a given distribution of z , the change-of-variable rule gives

$$p(X) = p(z) \left| \det \frac{\partial f(X)}{\partial X} \right|, \quad (2)$$

where $\frac{\partial f(X)}{\partial X}$ is the Jacobian matrix of f^{-1} at X . To obtain a complex distribution with expressive mapping, a series of invertible transformations are chained together. The relationship between input X and latent representation z becomes:

$$X \xleftrightarrow{f_1} h_1 \xleftrightarrow{f_2} h_2 \cdots \xleftrightarrow{f_K} h_K, \quad z \triangleq h_K \quad (3)$$

The sequence of invertible transformations f^{-1} in Equation 3 is known as a normalizing flow. We can generate samples X

by first sampling $z \sim p_\theta(z)$ and then computing $X = f^{-1}(z)$. Then, we can write the log-likelihood of X as

$$\log p_\theta(X) = \log p_\theta(z) + \log \left| \det \frac{\partial z}{\partial X} \right| \quad (4)$$

$$= \log p_\theta(z) + \sum_{i=1}^K \log \left| \det \frac{\partial h_i}{\partial h_{i-1}} \right| \quad (5)$$

To get a tractable determinant of the Jacobian matrix, the idea is to select the transformation whose Jacobian is a triangular matrix. The log-determinant is then simplified as

$$\log \left| \det \frac{\partial h_i}{\partial h_{i-1}} \right| = \operatorname{sum}(\log |\operatorname{diag}(\frac{\partial h_i}{\partial h_{i-1}})|). \quad (6)$$

Glow [16] is a normalizing flow-based generative model, which has achieved impressive performance for facial image synthesis. In Glow, each step of flow consists of three sub-steps: *actnorm*, *invertible 1×1 convolution*, and an *affine coupling layer*. Actnorm is an activation normalization layer that applies a scale and bias with data-dependent initialization. 1×1 convolution is a linear transformation layer for soft permutation. The affine coupling layer splits the input into two parts. One half of the input can be affinely transformed based the other half. The other half is kept unchanged, which leads to an easy reverse transformation.

MoGlow [1] extends Glow to skeleton sequences by describing the distribution of future poses in terms of an autoregressive model. It adds control signals to achieve control over the output and uses a recurrent neural network, an LSTM, to integrate information over time. The past and current control signal $C_{(t-T_h):(t)}$, and past human poses $\mathbb{X}_{(t-T_h):(t-1)}$ are conditioning information. MoGlow simply concatenates the current pose with additional conditioning information and feeds it into the LSTM in the affine coupling layers. The autoregressive model can be written as

$$p_\theta(\mathbb{X}|C) = p(\mathbb{X}_{(t_0-T_h):(t_0-1)} | C_{(t_0-T_h):(t_0-1)}) \cdot \prod_{t=t_0}^T p_\theta(X_t | \mathbb{X}_{(t-T_h):(t-1)}, C_{(t-T_h):(t)}, H_{t-1}), \quad (7)$$

where H_t is the latent state of the LSTM. In our proposed model, indicative features of the past poses are first extracted using spatial-temporal graph convolutional networks before being fed into the recurrent neural network, which will be explained in the following.

C. Graph Motion Glow

In MoGlow, the coordinates of the skeleton data are concatenated to one single feature vector per frame. With a spatial graph neural network (S-GCN), we instead convert the skeleton data into an undirected graph $G = (V, E)$, where V is the set of nodes and E is the set of edges, based on the natural spatial relationships between parts of the skeleton. Each joint marker is represented by a node and the bone between two connected markers by an edge. The skeleton graph used in this work is illustrated in Fig. 3.

Spatial graph convolutional networks extend convolutions, typically applied to images, to graph structures [43]. In the

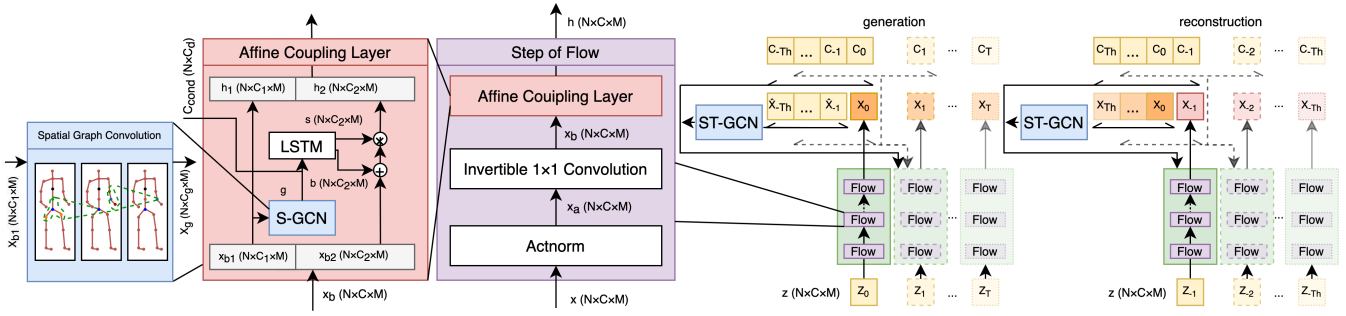


Fig. 2: Overview of the Graph Motion Glow for skeleton-based motion reconstruction and generation. The body markers are connected as a skeleton graph and fed into the step of flows the with control input and history input. The S-GCN in the affine transformation encodes the spatial relationships among markers, and the LSTM preserves the temporal information. The ST-GCN extracts features from past sequences.



Fig. 3: The spatial graph of a skeleton used in this work. Each graph node corresponds to the marker on the right [42]. The edges between nodes are defined based on natural joints relationships.

case of images, pixels within a rectangular neighborhood have a fixed spatial order, but no such order exists in a skeleton graph. To address this problem, Yan et al. proposed a spatial configuration partition strategy [27]. We follow the same idea to define a sampling function for convolution. The spatial configuration partition strategy divides each neighborhood of a node into three subsets, i.e., the neighbors located closer to the graph center, the neighbors that are farther away, and the node itself.

As shown in Fig. 3, we assume the chest node, marker-10, to be the center of the skeleton. For the neighbor set of marker-0 with kernel scale $D = 1$ (nodes within the green dotted line), there are three subsets: marker-9 (node closer to the center, green), marker-1 and marker-5 (nodes farther away, orange), and marker-0 (the node itself, blue). Similar neighborhood subsets are determined for each node in the skeleton graph for different kernel scales.

A spatial graph convolution can then be defined as:

$$y_i = \sum_{v_j \in S_i} \frac{x_j}{D_{v_i}(v_j)} w(l_{v_i}(v_j)), \quad (9)$$

where x_i is the feature vector of node v_i before the convolution, y_i is the corresponding vector after the convolution, S_i is the neighbor set of v_i , and w is a weight function that depends on the spatial partition strategy $l_{v_i}(v_j)$ that maps each node v_j to its corresponding subset. $D_{v_i}(v_j)$ is the number of nodes in the subset, which is used as a normalizing term to compute

an average per subset.

We also employ temporal connections that connect the same node in consecutive frames with temporal edges. For a past sequence, $v_{t,i}$, at time t , connects to $v_{t-1,i}$ and $v_{t+1,i}$ along the temporal dimension. The whole graph sequence is thus composed of a spatial graph and a temporal graph. As in ST-GCN [27], each layer of the network includes one S-GCN followed by a temporal convolution network (TCN) over the temporal domain, with residual connections.

An overview of the proposed graph motion glow framework is presented in Fig. 2, which includes each component transformation. Our flow model includes the three main reversible transformation layers in Glow, but extended to graph structures. We further use ST-GCN to extract features from the autoregressive history input. The input X and output z are represented as tensors of shape $[M \times C \times T_h]$ with spatial dimension M , channel dimension C and temporal dimension T_h . X_a and X_b denote intermediate results of the actnorm layer and invertible 1×1 convolution layer, which performs a soft permutation of channels. The affine coupling layer is more complex. The idea is to split the input channel-wise into two parts and transform one part based on the other part and conditioning information. Mathematically, we define the input X_b and output h of the affine coupling layer as concatenations $X_b = [X_{b1}, X_{b2}]$ and $h = [h_1, h_2]$. The coupling can then be written

$$[h_1, h_2] = [X_{b1}, (X_{b2} + \mathbf{b}) \odot \mathbf{s}], \quad (10)$$

where \odot is a Hadamard product and the scaling \mathbf{s} and bias \mathbf{b} terms are computed with S-GCN, ST-GCN and LSTM:

$$g_t = SGCN(X_{b1,t}), \quad (11)$$

$$p_t = STGCN(\hat{X}_{(t-T_h):(t-1)}), \quad (12)$$

$$[\mathbf{s}_t, \mathbf{b}_t] = LSTM(g_t, p_t, C_{(t-T_h):(t)}). \quad (13)$$

Here, S-GCN captures the spatial graph information g_t from markers in the current time step t , ST-GCN extracts spatial-temporal features p_t from the past frames of the sequence, and LSTM produces the scaling and bias with dependencies over time. At different steps of the flow, we use spatial graph convolutions with different graph kernel scales D to capture the hierarchical structure of the body.

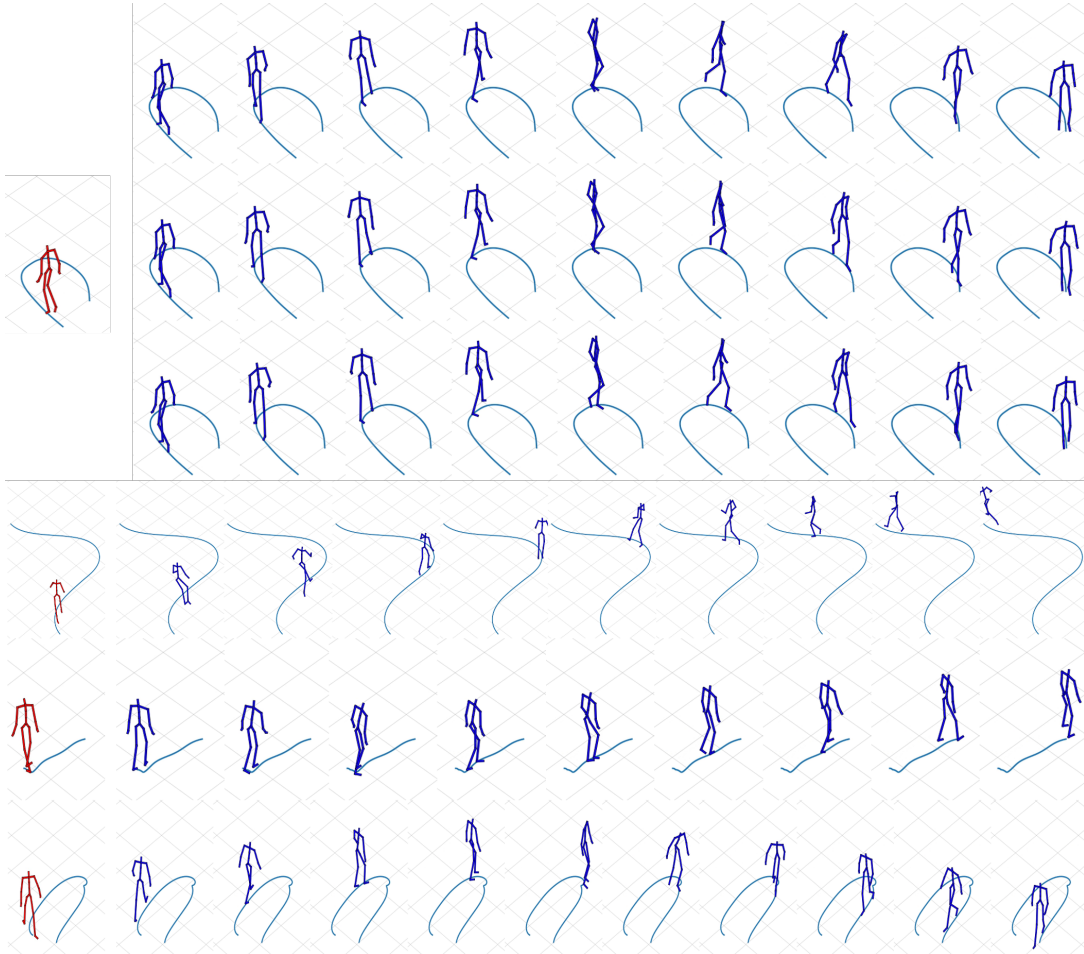


Fig. 4: Example sequences generated by STMG. The top three sequences are generated from the same past information. The bottom three sequences are generated given different past information.

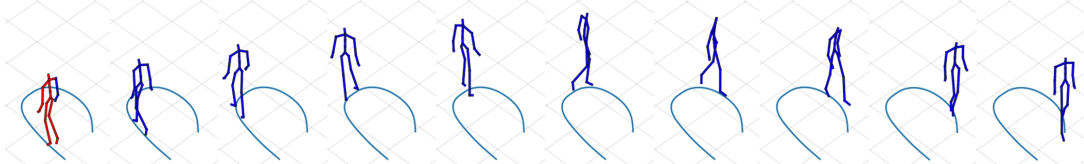


Fig. 5: Example sequence generated and reconstructed by STMG. In the past poses, the same as the first three rows of Fig. 4, the markers of the right arm are set to zero.

For generation, we exploit the fact that the flow model is reversible. Thus we can generate a new pose X_t from the trained graph motion glow model using a latent vector z_t (see Section III-B), the recent history of poses $\mathbb{X}_{(t-T_h):(t-1)}$ and control input $C_{(t-T_h):(t)}$ for the full sequence. The latent vector z_t is sampled independently from a standard Gaussian distribution. The generated X_t then becomes a part of the conditioning information for generating the next following pose X_{t+1} . During training, data was augmented by lateral mirroring and time-reversion. This allows an imperfect input with missing markers to be reconstructed with the same framework by reversing the generated sequences $\mathbb{X}_{(t_0):(T_h)}$ and control signal $C_{(t_0-T_h):(T_h)}$ to $\mathbb{X}_{(T_h):(t_0)}$ and $C_{(T_h):(t_0-T_h)}$. With these reversed sequences regarded as control information, an imperfect input can be reconstructed with generated markers now used to fill in the holes of the missing data.

IV. EXPERIMENTAL RESULTS

In this section, we evaluate the performance of our model. We first describe the dataset and experiment setup, and then continue with the qualitative and quantitative results on generating and reconstructing human motion samples of high fidelity. On basis of these, we discuss strengths and limitations of the baseline and proposed frameworks.

A. Dataset

We consider a human locomotion dataset preprocessed by [1], which pools samples from the Edinburgh Locomotion MOCAP[42], CMU Motion Capture [44], and HDM05 [45] datasets. The dataset is downsampled to 20 fps and sliced into fixed-length sequences of 80 frames with 50% overlap for training and augmented by lateral mirroring and time-reversal. Each data instance resembles a clip of a human

locomotion animation, i.e., movement including various gaits along different curved paths.

The motion is represented by the 3D coordinates of 21 joints at each frame, as displayed in Figure 3. In addition, there are 3 scalar control signals indicating forward, sideways and rotational velocities for each frame. To generate incomplete MoCap frames, we set some markers to zero with a binary matrix M_b in the past poses used for generation: $\hat{\mathbb{X}}_{(t_0-T_h):(t_0-1)} = M_b \odot \mathbb{X}_{(t_0-T_h):(t_0-1)}$. We preprocess these sequences with the following five settings: a) remove the markers on the right arm (markers 18, 19, 20), b) remove the markers on the left leg (markers 2, 3, 4), c) remove the markers on the right arm and left leg, d) randomly remove 4 markers, e) keep all markers intact.

B. Proposed Model and Ablations

The proposed graph motion glow model is trained with a 10-frame time window, similarly to MoGlow. In our experiments, all glow models were structured with 16 steps of flows. The structure of one step of flow is illustrated as the purple block in Fig. 2. Each affine coupling step consisted of a S-GCN layer and an LSTM with two layers, shown as the red block in Fig. 2. Note that ST-GCN is applied to the past history of frames with the same features used in each such step. The sizes of graph kernels for the S-GCN are 3, 5, and 7 for the first 10 flow steps, the following 4 steps, and the last 2 steps. We use a temporal kernel with a size of 9 and 512 hidden cells for each LSTM layer.

To assess the impact of design decisions, we trained three versions of the architecture on the human data. Each version has specific components disabled to examine their effects comparing with our full framework. We denote our proposed graph-based model, the spatial-temporal graph motion glow as "STMG". The first ablated configuration "SMG" only uses the spatial graph convolution networks without temporal convolution, i.e., the temporal convolution in the "ST-GCN" block in Fig. 2 is turned off. The second configuration uses no graph structure and is equivalent to the MoGlow baseline MoGlow, denoted as "MG".

C. Results and Discussions

We present examples of generated sequences in Fig. 4, which shows snapshots of every 10 frame from sequences of length 100. A video with generated examples can be found at this link⁴. The results demonstrate the diversity and quality of samples. The top three sequences are sampled from the same 10-frame seeding history but show different poses, which means the probabilistic model does not collapse into a stereotypical mode. The bottom three sequences are generated given different past poses and control inputs, which show our model incorporating spatial and temporal graph structure can generate long-term locomotion behavior of high fidelity and diversity. When the past information is incomplete, we further reconstructed the missing markers. In Fig. 5, we observe that the reconstructed markers (the right arm in the first frame of Fig. 5) fit the original skeleton well. We further evaluate the quality of generated and reconstructed motions

by performing both footsteps analysis (see [1] for detailed definitions) and bone-length analysis.

1) *Footstep Analysis*: Footsteps analysis is used to evaluate foot-sliding artifacts in locomotion synthesis. In footstep analysis, footsteps are detected as time intervals when the horizontal speed of the heel joints is below a tolerance value v_{tol} . Because of foot-sliding artifacts, the heels sometimes exhibit a spurious behavior with the horizontal speed exceeding a specified tolerance value. We incremented the tolerance v_{tol} in small steps. The total number of detected footsteps first rises, then reaches the maximum value. The number of identified footsteps decreases if the tolerance increasing further. Motion sequences with larger foot-sliding artifacts need a higher tolerance value to reach the maximum value of estimated footsteps.

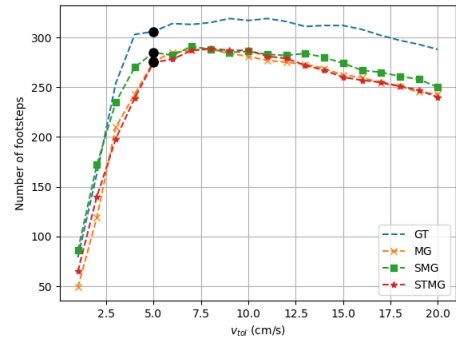


Fig. 6: Footstep analysis for samples generated with complete past input: footstep count f_{est} on tolerance value v_{tol} . Black dots indicate the location of v_{tol}^{95} .

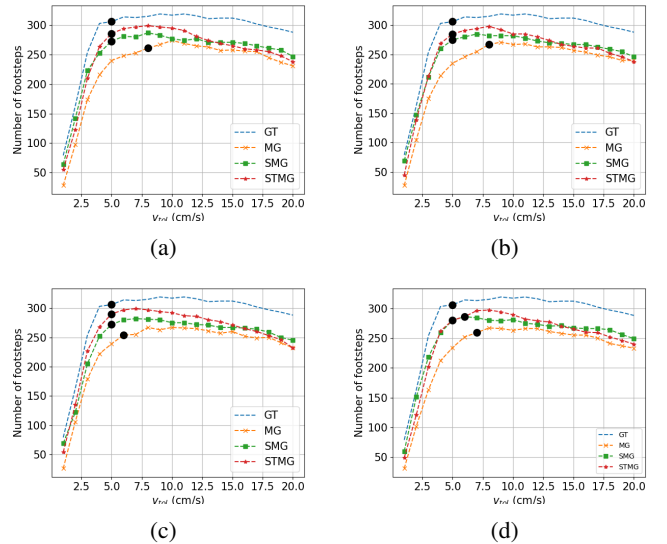


Fig. 7: Footstep analysis for incomplete past input: footstep count f_{est} for each tolerance v_{tol} . Black dots indicate v_{tol}^{95} . Graphs represent cases with missing markers (a) on the right arm (marker-18, 19, 20), (b) on the left leg (marker 2, 3, 4), (c) on the right arm and left leg, and (d) 4 randomly missing.

The estimated number of footsteps for different tolerance values in generated motions are shown in Fig. 6 and Fig.

⁴<https://kth.box.com/s/2vngw2tu1pg217cf9bo4s98fzobnccxb>

7. Each results is based on an average over 150 generated sequences. For evaluation, we detect the first tolerance value v_{tol}^{95} , for which at least 95% of the maximum number of footsteps are estimated. These values are shown as black dots in the figures. The total estimated number of footsteps, speed tolerance for capturing 95% steps, and the mean and standard deviation of the step duration are shown in Table I. We note that in Fig. 6, which assumes sequences without missing markers, the curves are close. The performance of the proposed method is comparable to the state-of-the-art baseline. However, when the given past data is incomplete, i.e., some markers are missing, we note in Fig. 7 that the curves of our proposed spatial temporal graph motion glow model are consistently closer to the curve of ground truth, illustrating the improvement in robustness of the proposed graph model. We also observe from the ablation study that STMG outperforms SMG, showing an extra performance gain from the temporal convolutions of the past frames.

| Miss | Model | f_{est} | v_{tol}^{95} | μ | σ |
|---------|-------|-----------|----------------|--------------|--------------|
| - | GT | 5 | 306 | 0.315 | 0.273 |
| - | MG | 5 | 276 | 0.298 | 0.318 |
| | SMG | 5 | 285 | 0.294 | 0.242 |
| | STMG | 5 | 275 | 0.316 | 0.267 |
| RA | MG | 8 | 261 | 0.380 | 0.283 |
| | SMG | 5 | 273 | 0.281 | 0.231 |
| | STMG | 5 | 286 | 0.302 | 0.249 |
| LL | MG | 8 | 267 | 0.382 | 0.362 |
| | SMG | 5 | 275 | 0.287 | 0.248 |
| | STMG | 5 | 285 | 0.314 | 0.275 |
| RA & LL | MG | 6 | 254 | 0.327 | 0.318 |
| | SMG | 5 | 272 | 0.275 | 0.217 |
| LL | STMG | 5 | 290 | 0.306 | 0.267 |
| | MG | 7 | 256 | 0.357 | 0.329 |
| | SMG | 5 | 280 | 0.307 | 0.252 |
| R4M | STMG | 6 | 286 | 0.315 | 0.253 |

TABLE I: Results of foot-step analysis for motion generation: total number of footsteps f_{est} , speed tolerance for capturing 95% steps v_{tol}^{95} , mean μ and standard deviation σ of step-duration. We remove a few markers in the past poses. RA: right arm; LL: left leg; R4M: random 4 markers. The numbers closest to the ground truth are shown in bold.

| Miss | Model | Generation | | Reconstruction | |
|---------|-------|--------------|--------------|----------------|--------------|
| | | RMSE | σ | RMSE | σ |
| - | MG | 0.597 | 0.067 | - | - |
| | SMG | 0.191 | 0.039 | - | - |
| | STMG | 0.779 | 0.073 | - | - |
| RA | MG | 279301 | 7.532 | 137325 | 4.471 |
| | SMG | 0.842 | 0.051 | 2.589 | 0.075 |
| | STMG | 0.881 | 0.079 | 0.890 | 0.059 |
| LL | MG | 834081 | 12.558 | 690231 | 7.663 |
| | SMG | 1.935 | 0.050 | 3.701 | 0.130 |
| | STMG | 0.909 | 0.080 | 2.336 | 0.122 |
| RA & LL | MG | 537036 | 8.788 | 2404434 | 7.533 |
| | SMG | 2.134 | 0.052 | 3.687 | 0.103 |
| | STMG | 0.917 | 0.080 | 1.864 | 0.093 |
| R4M | MG | 787638 | 12.138 | 80931 | 2.897 |
| | SMG | 0.542 | 0.044 | 6.589 | 0.092 |
| | STMG | 0.938 | 0.080 | 1.072 | 0.065 |

TABLE II: Results of bone-length analysis for human motion generation and reconstruction. The best values are in bold.

2) *Bone-length Analysis*: The human skeleton is represented by joint coordinates, without explicit consideration of

spatial constraints such as the length of bones. To evaluate the data quality on this aspect, [1] performs bone-length analysis to detect artifacts such as flying-apart joints. The analysis looks at the bone-length Root Mean Squared Error (RMSE) bl_{rmse} (cm) and standard deviation bl_{σ} (cm^2), which represent the consistency, stability and quality of generation and reconstruction. Table II reports the results on the bone-length analysis. Given complete past frames, all competing models achieve relative small RMSE and σ , with SMG outperforming the other two.

Again, the performance of baseline and our approaches differs significantly when the models are tasked to generalize from untrained imperfect data. MG performs poorly on this indicator, exhibiting huge bone-length artifacts. In contrast, for both SMG and STMG, the RMSE and σ of bone-lengths are still small.

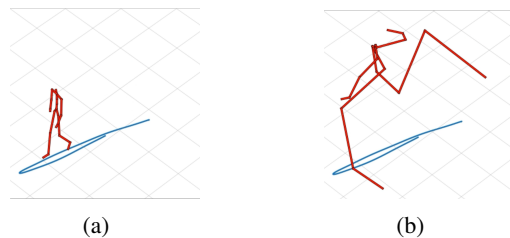


Fig. 8: Examples of pose generated by MG. (a) with complete past poses. (b) with incomplete past poses.

A closer look at the samples reveals unstable generation and reconstruction with MG. A typical example is shown in Fig. 8. When the conditioning information includes complete past poses, the generated pose is stable. However, when the conditioning information is incomplete, with four random markers missing, we observe that some joints fly apart. This situation occurs occasionally. With incomplete conditioning information, there exists a distribution shift that may result in bone-stretching or joints flying apart, since the feed-forward neural networks are sensitive to the shifts of the data domain. For SMG and STMG, the generated poses are always stable and the RMSE and σ of bone lengths are small, since we exploit the correlation of human skeletons to enforce the bias with graph structures.

V. CONCLUSION

In this paper, we propose a graph-based normalizing flow model to tackle the problem of human motion generation and reconstruction. This new modelling framework has the following main advantages: (1) It is an extension of MoGlow that is probabilistic and allows inference of the exact likelihood. (2) It utilizes spatial-temporal graph convolutional networks to improve the robustness of generation. To the authors' knowledge, this is the first work where the graph-based normalizing flow is used to generate and reconstruct human motion. From the results it can be concluded that it overcomes some of the limitations in earlier models. In the future, we plan to extend the graph-based motion glow model to multiple scales to tackle more complex motions.

ACKNOWLEDGEMENTS

This research has received funding from the EC Horizon 2020 research and innovation program under grant agreement n. 824160 (EnTimeMent).

REFERENCES

- [1] G. E. Henter, S. Alexanderson, and J. Beskow, "Moglow: Probabilistic and controllable motion synthesis using normalising flows," *ACM Transactions on Graphics (TOG)*, vol. 39, no. 6, pp. 1–14, 2020.
- [2] R. Murakami, L. Y. Morales Saiki, S. Satake, T. Kanda, and H. Ishiguro, "Destination unknown: Walking side-by-side without knowing the goal," in *Proceedings of the 2014 ACM/IEEE International Conference on Human-Robot Interaction*, ser. HRI '14. New York, NY, USA: Association for Computing Machinery, 2014, p. 471–478.
- [3] K. Fragkiadaki, S. Levine, P. Felsen, and J. Malik, "Recurrent network models for human dynamics," in *Proceedings of the IEEE International Conference on Computer Vision*, 2015, pp. 4346–4354.
- [4] J. Butepage, M. J. Black, D. Kragic, and H. Kjellström, "Deep representation learning for human motion prediction and classification," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2017, pp. 6158–6166.
- [5] J. Martinez, M. J. Black, and J. Romero, "On human motion prediction using recurrent neural networks," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2017.
- [6] Z. Wang, J. Chai, and S. Xia, "Combining recurrent neural networks and adversarial training for human motion synthesis and control," *IEEE transactions on visualization and computer graphics*, vol. 27, no. 1, pp. 14–28, 2019.
- [7] H. Ahn, J. Kim, K. Kim, and S. Oh, "Generative autoregressive networks for 3d dancing move synthesis from music," *IEEE Robotics and Automation Letters*, vol. 5, no. 2, pp. 3500–3507, 2020.
- [8] I. Habibie, D. Holden, J. Schwarz, J. Yearsley, and T. Komura, "A recurrent variational autoencoder for human motion synthesis," in *28th British Machine Vision Conference*, 2017.
- [9] T. Kucherenko, J. Beskow, and H. Kjellström, "A neural network approach to missing marker reconstruction in human motion capture," *arXiv preprint arXiv:1803.02665*, 2018.
- [10] Z. Li, Y. Zhou, S. Xiao, C. He, Z. Huang, and H. Li, "Auto-conditioned recurrent networks for extended complex human motion synthesis," *arXiv preprint arXiv:1707.05363*, 2017.
- [11] S. Yan, Z. Li, Y. Xiong, H. Yan, and D. Lin, "Convolutional sequence generation for skeleton-based action synthesis," in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2019.
- [12] N. Sadoughi and C. Busso, "Novel realizations of speech-driven head movements with generative adversarial networks," in *2018 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 2018, pp. 6169–6173.
- [13] T. Kucherenko, P. Jonell, S. van Waveren, G. E. Henter, S. Alexandersson, I. Leite, and H. Kjellström, "Gesticulator: A framework for semantically-aware speech-driven gesture generation," in *Proceedings of the 2020 International Conference on Multimodal Interaction*, 2020.
- [14] S. Ghorbani, C. Wloka, A. Etemad, M. A. Brubaker, and N. F. Troje, "Probabilistic character motion synthesis using a hierarchical deep latent variable model," in *Computer Graphics Forum*, vol. 39, no. 8. Wiley Online Library, 2020, pp. 225–239.
- [15] D. Greenwood, S. Laycock, and I. Matthews, "Predicting head pose from speech with a conditional variational autoencoder." ISCA, 2017.
- [16] D. P. Kingma and P. Dhariwal, "Glow: Generative flow with invertible 1x1 convolutions," in *Advances in neural information processing systems*, 2018, pp. 10215–10224.
- [17] S. Alexanderson, G. E. Henter, T. Kucherenko, and J. Beskow, "Style-controllable speech-driven gesture synthesis using normalising flows," in *Computer Graphics Forum*, vol. 39, no. 2, 2020, pp. 487–496.
- [18] S.-J. Peng, G.-F. He, X. Liu, and H.-Z. Wang, "Hierarchical block-based incomplete human mocap data recovery using adaptive nonnegative matrix factorization," *Computers & Graphics*, vol. 49, 2015.
- [19] Z. Wang, S. Liu, R. Qian, T. Jiang, X. Yang, and J. J. Zhang, "Human motion data refinement unitizing structural sparsity and spatial-temporal information," in *2016 IEEE 13th International Conference on Signal Processing (ICSP)*. IEEE, 2016, pp. 975–982.
- [20] Ø. Gløersen and P. Federolf, "Predicting missing marker trajectories in human motion data using marker intercorrelations," *PLoS one*, vol. 11, no. 3, p. e0152616, 2016.
- [21] C.-H. Tan, J. Hou, and L.-P. Chau, "Human motion capture data recovery using trajectory-based matrix completion," *Electronics Letters*, vol. 49, no. 12, pp. 752–754, 2013.
- [22] W. Hu, Z. Wang, S. Liu, X. Yang, G. Yu, and J. J. Zhang, "Motion capture data completion via truncated nuclear norm regularization," *IEEE Signal Processing Letters*, vol. 25, no. 2, pp. 258–262, 2017.
- [23] Q. Cui, B. Chen, and H. Sun, "Nonlocal low-rank regularization for human motion recovery based on similarity analysis," *Information Sciences*, vol. 493, pp. 57–74, 2019.
- [24] U. Mall, G. R. Lal, S. Chaudhuri, and P. Chaudhuri, "A deep recurrent framework for cleaning motion capture data," *arXiv preprint arXiv:1712.03380*, 2017.
- [25] Q. Cui, H. Sun, Y. Li, and Y. Kong, "A deep bi-directional attention network for human motion recovery," in *IJCAI*, 2019, pp. 701–707.
- [26] S. Lohit, R. Anirudh, and P. Turaga, "Recovering trajectories of unmarked joints in 3d human actions using latent space optimization," in *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*, 2021, pp. 2342–2351.
- [27] S. Yan, Y. Xiong, and D. Lin, "Spatial temporal graph convolutional networks for skeleton-based action recognition," in *Proceedings of the AAAI conference on artificial intelligence*, vol. 32, no. 1, 2018.
- [28] M. Li, S. Chen, X. Chen, Y. Zhang, Y. Wang, and Q. Tian, "Action-structural graph convolutional networks for skeleton-based action recognition," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2019, pp. 3595–3603.
- [29] K. Cheng, Y. Zhang, X. He, W. Chen, J. Cheng, and H. Lu, "Skeleton-based action recognition with shift graph convolutional network," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2020, pp. 183–192.
- [30] F. Yang, W. Yin, T. Inamura, M. Björkman, and C. Peters, "Group behavior recognition using attention-and graph-based neural networks," in *Proceedings of the 24th European Conference on Artificial Intelligence*, 2020.
- [31] C. Si, W. Chen, W. Wang, L. Wang, and T. Tan, "An attention enhanced graph convolutional lstm network for skeleton-based action recognition," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2019, pp. 1227–1236.
- [32] K. Thakkar and P. Narayanan, "Part-based graph convolutional network for action recognition," *arXiv preprint arXiv:1809.04983*, 2018.
- [33] Y.-H. Wen, L. Gao, H. Fu, F.-L. Zhang, and S. Xia, "Graph cnns with motif and variable temporal block for skeleton-based action recognition," in *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 33, no. 01, 2019, pp. 8989–8996.
- [34] M. Li, S. Chen, Y. Zhao, Y. Zhang, Y. Wang, and Q. Tian, "Dynamic multiscale graph neural networks for 3d skeleton based human motion prediction," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2020, pp. 214–223.
- [35] A. Jain, A. R. Zamir, S. Savarese, and A. Saxena, "Structural-rnn: Deep learning on spatio-temporal graphs," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2016.
- [36] L. Zhao, X. Peng, Y. Tian, M. Kapadia, and D. N. Metaxas, "Semantic graph convolutional networks for 3d human pose regression," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2019, pp. 3425–3435.
- [37] F. Yang, W. Yin, M. Björkman, and C. Peters, "Impact of trajectory generation methods on viewer perception of robot approaching group behaviors," in *2020 29th IEEE International Conference on Robot and Human Interactive Communication (RO-MAN)*, 2020, pp. 509–516.
- [38] Y. Zhao and Y. Dou, "Pose-forecasting aided human video prediction with graph convolutional networks," *IEEE Access*, vol. 8, 2020.
- [39] J. Liu, A. Kumar, J. Ba, J. Kiro, and K. Swersky, "Graph normalizing flows," *arXiv preprint arXiv:1905.13177*, 2019.
- [40] L. Dinh, D. Krueger, and Y. Bengio, "Nice: Non-linear independent components estimation," *arXiv preprint arXiv:1410.8516*, 2014.
- [41] L. Dinh, J. Sohl-Dickstein, and S. Bengio, "Density estimation using real nvp," *arXiv preprint arXiv:1605.08803*, 2016.
- [42] I. Habibie, D. Holden, J. Schwarz, J. Yearsley, and T. Komura, "A recurrent variational autoencoder for human motion synthesis," in *Proceedings of the British Machine Vision Conference (BMVC)*, 2017.
- [43] J. Dai, H. Qi, Y. Xiong, Y. Li, G. Zhang, H. Hu, and Y. Wei, "Deformable convolutional networks," in *Proceedings of the IEEE international conference on computer vision*, 2017, pp. 764–773.
- [44] [Online]. Available: <http://mocap.cs.cmu.edu/>
- [45] M. Müller, A. Baak, and H.-P. Seidel, "Efficient and robust annotation of motion capture data," in *Proceedings of the ACM SIGGRAPH/Eurographics Symposium on Computer Animation (SCA)*.