

# Learn to Predict How Humans Manipulate Large-sized Objects from Interactive Motions

Weilin Wan<sup>1</sup>, Lei Yang<sup>1,4</sup>, Lingjie Liu<sup>2</sup>, Zhuoying Zhang<sup>1</sup>, Ruixing Jia<sup>1</sup>, Yi-King Choi<sup>1,4</sup>, Jia Pan<sup>1</sup>, Christian Theobalt<sup>2</sup>, Taku Komura<sup>1</sup> and Wenping Wang<sup>3</sup>

**Abstract**—Understanding human intentions during interactions has been a long-lasting theme, that has applications in human-robot interaction, virtual reality and surveillance. In this study, we focus on full-body human interactions with large-sized daily objects and aim to predict the future states of objects and humans given a sequential observation of human-object interaction. As there is no such dataset dedicated to full-body human interactions with large-sized daily objects, we collected a large-scale dataset containing thousands of interactions for training and evaluation purposes. We also observe that an object’s intrinsic physical properties are useful for the object motion prediction, and thus design a set of object dynamic descriptors to encode such intrinsic properties. We treat the object dynamic descriptors as a new modality and propose a graph neural network, HO-GCN, to fuse motion data and dynamic descriptors for the prediction task. We show the proposed network that consumes dynamic descriptors can achieve state-of-the-art prediction results and help the network better generalize to unseen objects. We also demonstrate the predicted results are useful for human-robot collaborations.

**Index Terms**—Intention Recognition, Human-Robot Collaboration, Datasets for Human Motion

## I. INTRODUCTION

Predicting human intentions to manipulate or carry objects, or more specifically, the future states of such interactions from a given sequential observation, has applications in robotic prosthesis and human-robot interaction [5], [35], [20], [19], [18], where the accuracy of such prediction can strongly affect the embodied perception and VR/3D gaming [15], where compensation of latency through such prediction is needed for improving the user experience. Despite the high demand, research for such prediction has been limited to small-sized objects [32], [6] or complex scenes [4], [39], [27]. Research for interactions with larger objects, such as chairs or boxes, that involve full-body motions, is rather under-explored and requires considering physics for reliable prediction.

Learning the physical interactions between the human body and arbitrary objects is hard due to the large variation of

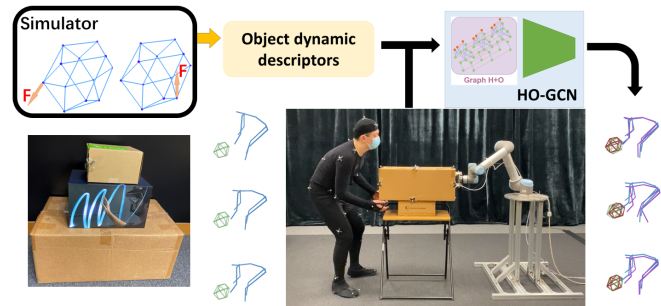


Fig. 1. We introduce a large-scale dataset on full-body human interactions with large-sized objects, such as chairs and boxes, using a motion capture system. Our goal is to predict the human-object motion at future time steps (where predicted human and object are in purple and red, respectively). To this end, we propose to leverage the *object dynamic descriptors* and design a neural network, *HO-GCN*, to fuse data of different modalities. We also showcase a human-robot collaborative task to validate the proposed method.

object geometries, the complex physics involved, and a large amount of training needed [11], [2], [14], [10]. Ehsani et al. [10] propose to integrate physics into a neural network and leverage it for object motion prediction. Nevertheless, the exact geometry of the object is required for the neural framework. They also require conducting physical simulation within the training loop, which is expensive and time-consuming for good generalization.

We make two efforts to address the aforementioned challenges. First, in response to the lack of datasets concerning interactions between humans and large-sized objects, we construct such a dataset. It contains 384K frames from 508 full-body motion capture (MoCap) videos, involving 12 daily objects from 6 categories and different actions. Not only included are the human body motions but also the 6 degrees-of-freedom (DOF) motions of the objects that are represented by 12 keypoints. The interaction take place in diverse contexts, such as transporting the objects forward or backward, rotating the objects, or more complex compositional motions (e.g., lift and carry forward), depending on the affordances and functionality of the object. The data were fully preprocessed and will be available to the community upon the publication of the paper. To the best of our knowledge, this is the first large-scale dataset focusing full-body interactions with large-sized objects.

To avoid frequent calls for physical simulators during training, we propose a novel descriptor that encodes the object’s intrinsic dynamics obtained from simulations. In particular, a conceptual model for each category of objects (e.g., chairs)

Manuscript received: September 9, 2021; Accepted: November 22, 2021.

This paper was recommended for publication by Editor Angelika Peer upon evaluation of the Associate Editor and Reviewers’ comments.

<sup>1</sup> Department of Computer Science, The University of Hong Kong, Hong Kong SAR

<sup>2</sup> Max-Planck-Institute for Informatics, Saarbruecken, Germany

C. Theobalt was supported by the ERC Consolidator Grant 4DRepLy (770784). L. Liu was supported by Lise Meitner Postdoctoral Fellowship.

<sup>3</sup> Department of Computer Science and Engineering, Texas A&M University, TX, USA

<sup>4</sup> Centre for Garment Production Limited, Hong Kong SAR

This work was partially supported by the Innovation and Technology Commission of the HKSAR Government under the InnoHK initiative.

Digital Object Identifier (DOI): see top of this page.

is constructed, which is represented by a set of keypoints abstracting the general shape of this object category. Then, we simulate its responses under forces as a rigid-body system and acquire the intrinsic dynamic properties. Thus, the *object dynamic descriptors* are defined as how the keypoint-based conceptual model is transformed under given forces. As shown in our experiments, our method using the object dynamic descriptors can achieve state-of-the-art performances regarding the prediction of object motions even when the object instance is unseen in the training set. For human-robot collaborative tasks like lifting a box or handing over an object, this is crucial as reliable prediction of the object motions can enable the robot to provide desirable assistive operations to the human partners.

We also design a novel graph convolutional neural network (HO-GCN) that takes as input the human-object interactive motions as well as the object dynamic descriptors as shown in Fig. 1. We adopt the spatial-temporal graph convolution proposed in [37] to learn motion features and predict the future interactive motions based on the input sequence. We evaluate the proposed method as well as several state-of-the-art methods on our collected dataset. We also showcase that the predicted interactive motion can enable robot assistance in labor-intensive tasks, such as transporting a box.

Our technical contributions are threefold:

- 1) We contribute the *first* large-scale dataset concerning human full-body interactions with large-sized daily objects;
- 2) We propose to consider the object’s intrinsic dynamics as an extra modality for enhancing prediction of the object future poses during human-object interactions;
- 3) We design a novel graph convolutional neural network that fuses the observed human-object interaction sequence and the object intrinsic dynamics for the prediction task.

## II. RELATED WORK

### A. Learning from human skeleton data

Recognizing human activities from skeletal data receives increasing research attention as the acquisition of human skeletons becomes easier, such as employing motion capture techniques. With the development of deep learning, many attempts are made to learn features from the temporal sequences of the human skeleton data for action recognition or motion prediction. These works can be roughly categorized into RNN/LSTM-based methods which recursively process the temporal information [31], [38], [21], [6], and CNN-based methods which maps the temporal information into hidden features using convolution [33], [9], [22], [8], [3].

Since a human skeleton is naturally a graph structure, Graph Convolutional Networks (GCNs) are proposed for processing human skeleton data as well. Yan et al. [37] propose spatial-temporal GCN to efficiently extract the spatial and temporal features from the skeleton inputs for action recognition. Li et al. [23] extend the human skeleton graph in GCN with extra links to capture more dependencies and explore significant features of movements. Even GCNs have shown to be effective in representing human skeletons as graphs, only a few attempts

have been made to couple human skeletons and objects in GCNs. Kim et al. [16] incorporate a single point containing objects’ positional information to the human skeleton graph for classifying human actions. Cui et al. [7] propose to learn additional connectivity among joints besides their natural linkages for motion prediction.

In our work, we design HO-GCN, a novel architecture that fuses object information and human skeleton information to predict the 6DOF motion of the object. We compare our proposed method with C-TE [3] and CAHMP [6]. The former is a convolution-based method using a mirrored encoder-decoder framework to process the temporal inputs and generate predicted skeleton motion sequences. The latter is a recently proposed RNN-based method for predicting interactive motions between humans and the static environments or smaller objects.

### B. Human-object interaction dataset

Existing 3D human datasets are mainly designed for action recognition or human motion prediction tasks [13], [25], [29], [36], [28], which contain a limited number of human-object interaction cases and do not provide 3D ground truth position of objects. There have also been works focusing on hand-object interactions. Although these datasets usually provide sufficient spatial information of the objects, they focus on hand-held objects that do not involve full-body movements in the sequences. Ehsani et al. [10] propose a dataset with 174 RGB hand-object interaction videos and 6D object pose annotation. Grasping Actions with Body (GRAB) [32] is a large-scale dataset focusing on human grasping actions, which presents 1334 videos with body information and object 3D models. First-Person Hand Action (FPHA) [12] provides 1175 RGB videos with hand positions and 6D object poses annotated.

Whole-Body Human Motion Database [27], [26] includes a number of samples including human interactions with the scene context involving tables, cups, ladders, etc. However, the number of direct human operations on these large objects (e.g., the table) is limited. This motivates us to propose a dataset for understanding full-body interactions with large-sized daily objects such as tables, chairs, and boxes. In these scenarios, human poses need to be adapted to the object properties, such as shape or center of mass, making it different from previous works that only model interactions with small or thin objects (e.g., golf clubs in [24]). We also anticipate increasing attention from the community to use such a dataset for learning a prior of human-object joint motions or predicting motions for human-robot interactions.

### C. Predicting human intention in human-robot cooperation

A bulk of literature has contributed to predicting human intentions or motions for efficient and safe human-robot collaborations. An early work [30] proposes a method based on the Gaussian mixed model to learn from human demonstrations and thus adjust the robot’s role in the human-robot table lifting task. To predict intentions, Hidden Markov Models are widely used (e.g., [34], [20] for recognizing human movements

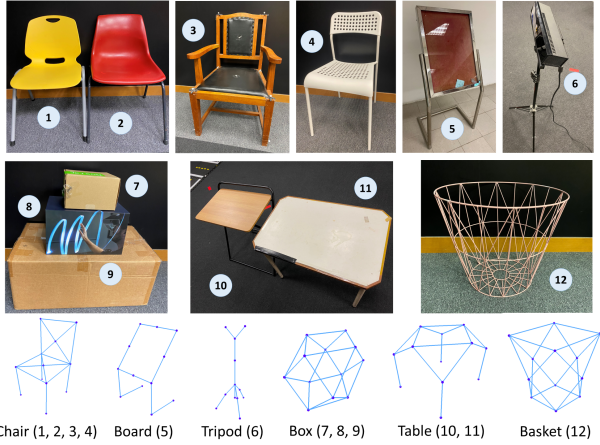


Fig. 2. The object instances in our dataset are shown in the top two rows with Object IDs next to them. The conceptual models used in the simulation are shown in the bottom row. Each conceptual model corresponds to an object class. They are modeled as rigid-body systems, and the drawn linkage is for illustration purposes only.

or interaction selection). Zhao et al. [40] proposed a recurrent network for imitation learning to accomplish object handover tasks between a human and a robot. In this paper, we focus on predicting interactive motions as a human manipulating a large-sized object, and showcase that the output of the proposed prediction method can be leveraged for human-robot collaboration.

### III. METHODOLOGY

#### A. Notations and problem formulation

The inputs to our prediction task are (1) a sequence of 3D human skeletons  $\{\mathbf{X}_t\}_{t=0}^K$ , (2) a sequence of 3D positional information of the object keypoints  $\{\mathbf{P}_t\}_{t=0}^K$ , and (3) the object dynamic descriptors  $\mathbf{O}$ . The expected outputs are the 6DoF pose changes of the object  $\{\Delta_t\}_{t=0}^K$  and the human skeleton motion  $\{\mathbf{X}_t\}_{t=0}^K$ . In what follows, we first detail the definition and acquisition of the object dynamic descriptors  $\mathbf{O}$  defined for objects and then present the network design.

#### B. Object representation and dynamic descriptors

**Geometric representation.** As many objects from the same class share high similarity in the overall geometry, we assume that objects from the same class can be abstracted by the same geometric conceptual model, and share common object dynamics. The conceptualized model is defined as a set of sparse keypoints that delineate the geometry of the object,  $\mathbf{P} = \{\mathbf{p}_m \in \mathbb{R}^3\}_{m=1}^M$ . The conceptualized models of all object categories from the dataset are shown in the bottom row of Fig. 2. The object dynamic descriptors are defined on the set of the keypoints of the conceptual model. We assume, for convenience, all conceptual models have the same number of keypoints and thus the same dimensionality of the dynamic descriptors. This allows training on all objects we collected in our dataset.

**Object dynamic descriptors.** Object dynamics are a set of intrinsic physical properties that reflect the resulting object

motions caused by external forces. In this study, we use  $\Delta = g(F)$  to describe how the given unit force  $F$  will cause the pose change  $\Delta$  of the (conceptual) object. To reduce the sampling space for the shape and the force, we consider the object as a rigid-body system. Thus, we can perform physical simulations to obtain the relationship between object motions and the forces. Since forces form a vector space and are linear to their resultant motion, a set of suitable bases of the force space shall suffice to describe this intrinsic dynamic relationship.

We define a set of unit forces along a candidate set  $\mathbb{D}$  of discrete directions, i.e. forward, backward, left, right and up, in the local frame of the object. A force applied to keypoint  $\mathbf{p}_m$  is denoted  $F_m^j$  where  $j$  indicates a direction from the candidate set. Then, we collect the resulting pose changes  $\Delta_m^j$  subject to the applied force  $F_m^j$  from simulations to describe the object dynamics. During each simulation, the force is applied to the object within an infinitesimal time interval and the object is modeled as a rigid-body system with a pre-defined floor constraint and the gravity. Finally, the proposed dynamic descriptor for an object is obtained as  $U = \{\Delta_m^j | m = 1 \dots M, j \in \mathbb{D}\}$  so that it can be fed to common neural networks that consume this semi-structured, discrete information as input.

#### C. Network structure

To achieve the described task, we propose a novel human-object graph convolutional neural network (HO-GCN) which consists of five major branches as shown in Figure 3. The *human branch* (top left) is designed for local feature learning from the spatial-temporal graph  $H$ . The spatial dimension of graph  $H$  is a human skeleton that has  $N$  nodes representing human body parts. The temporal dimension of the graph is formed by connecting each node of a human skeleton to its counterpart in temporally adjacent skeletons. Specifically, this branch processes human skeleton data by applying to  $H$  the spatial-temporal graph convolutions (denoted *STGConv*) using the distance partitioning strategy [37].

The *object pre-processing branch* aims to fuse the dynamic descriptors and the object keypoints. Specifically, we represent an object as two parts: 1) a  $D$ -dim vector  $U$  that represents the dynamic descriptors and 2) 3D positional differences between each object keypoint and both hands in the corresponding input frames. For better mining the relations among keypoint positions and descriptors, we broadcast and concatenate the keypoint coordinates with the dynamic descriptors and further process them using a convolution-based object encoder to obtain the object motion features.

The central *human-object fusion branch* then takes the concatenation of human and object motion features from previous branches, and applies spatial-temporal graph convolutions to them based on graph  $H+O$ . In this graph structure, we connect the object keypoints to three joints: left and right hands as well as the hip of a human body, as the movements of the hand and hip joints provides prominent hints related to human-object contact and the proximity between the human and the object for predicting the human's intention. This branch then outputs

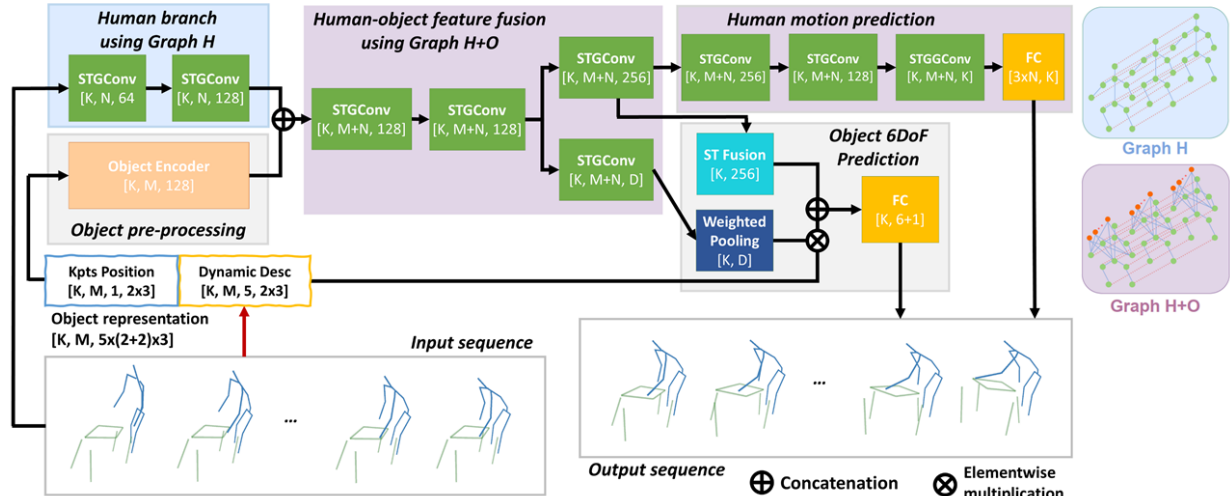


Fig. 3. The architecture of HO-GCN. The network consists of five major blocks. The two parallel blocks at the left end pre-process the human skeleton motion and the object information, and extract respective motion features. The central block using Graph  $H+O$  fuses the extracted features from the human and object, and performs graph convolutions to achieve global understanding of the input sequence. The upper and lower blocks at the right side respectively predict the human skeletons and the 6DoF changes of the objects in the future frames. *STGConv* and *FC* refers to spatial-temporal graph convolution and fully-connected layers, respectively.  $D$  equals  $M \times 5 \times 6$  in our case, where there are 5 candidates of discrete directions.

two tensors for predicting the human and object motions, respectively.

For predicting the 6DoF pose changes of the object in the future frames, we feed the output tensor with  $D$  channels from the central  $H+O$  branch to the *weighted pooling* module (as shown in Figure 3). This module is designed for scaling the dynamic descriptors, which are derived by applying unit forces to objects in a simulator. It maps the input tensor to the weighting factors which are then multiplied with the dynamic descriptor vector.

We also employ a spatial-temporal fusion module (*ST Fusion*) to recombine the respective information in spatial and temporal dimensions. This module yields an ST-fused feature containing information about the human-object interaction and is concatenated with the weighted dynamic descriptor vector. The concatenated features are then fed into a fully-connected (FC) layer for 6DoF regression. We additionally require the FC layer to regress a probability score  $c$  indicating if the object is in motion ( $c = 1$ ) and vice versa ( $c = 0$ ).

The *human motion prediction branch* performs a series of spatial-temporal graph convolutions based on graph  $H+O$  and eventually reduces the channel size of the graph convolution tensor to the number of  $K$  output frames. Finally, each of the  $K$  channels of the tensor is transformed (via a FC layer) to generate  $K$  human poses for  $K$  output frames, respectively.

**Loss function.** The loss function is formulated as:

$$L = \lambda_1 \left( \sum_t \|\hat{c}_t - c_t\|_2 + \|\hat{\Delta}_t - \Delta_t\|_2 + \sum_{p_i \in \mathbf{P}_t} \|\hat{p}_i - p_i\|_2 \right) + \lambda_2 \sum_{x_i \in \mathbf{X}_t} \|\hat{x}_i - x_i\|_2, \quad (1)$$

where notations with a hat, e.g.,  $\hat{x}$ , denote the predicted value, and those without a hat stand for the ground-truth;  $c = \{1, 0\}$  indicates the object in motion or not;  $\Delta$  is the pose change;  $p_i$

Translation
Push forward; Pull backward; Move to the left / right (154)
Rotation
Out-of-plane: Tilt to left / right / front / back sides (90)
In-plane: Rotate clockwise / anticlockwise (90)
Compositional actions
Lift and carry Forward (71)
Lift and place to left/right (60)
Lift and rotate up/down (43)

TABLE I  
HUMAN-OBJECT INTERACTIVE MOTIONS IN OUR COLLECTED DATASET.  
NUMBERS IN PARENTHESIS INDICATE THE NUMBER OF VIDEO SAMPLES.

and  $x_i$  are the object keypoints and human body joints. We set the balance coefficients  $\lambda_{i=\{1,2\}}$  to 1.0 and 0.5, respectively, for training our network.

#### IV. HUMAN-OBJECT INTERACTION DATASET

**Data collection.** We collected a dataset of 508 human-object interaction videos, 384K frames in total recorded by OptiTrack [1], a motion capture system. The average length of the recorded motions in our dataset is 6.3s with a standard deviation of 1.4s. Observing how humans may interact the large-size objects in the daily life, we design a set of interactive actions between humans and objects; see Tab. I. Six actors of different shapes, including two females and four males, participated in data collection process.

Fig. 2 shows the 12 large-size objects frequently seen in daily living, i.e., four chairs, a standing board, a tripod, three boxes of different sizes, two tables and a basket. Each of the objects is represented by a geometric abstraction of 12 keypoints delineating its shape, also shown in Fig. 2. We found 12 keypoints are sufficient to describe the geometry of the objects without too much computational overhead or memory consumption. We follow the keypoint configuration



of the conceptual model to attach twelve reflective markers on the real-world objects for tracking their rigid motions.

**Data processing.** We used the OptiTrack motion capture system (at a frame-rate of 120 Hz) to track the human-object interactions. Each MoCap video recorded a complete process of an actor performing given actions and captured the motion data of the human body skeleton and the pose of the target object at each frame.

We use a sliding window of 240 frames and a step size of 12 frames to extract sequences from a given MoCap video. Thus, each extracted sequence has 20 frames equivalent to a 2-sec motion. We use the first ten frames (1 sec) as input and predict human-object interaction in the last ten frames (the other 1 sec). We always ensure that the object at the  $K$ -th frame ( $K = 10$ ) in the input sequence to have a small motion within a threshold. For prediction, we label each frame in the extracted sequence with the 6DOF pose change  $\Delta_{gt}$  of the object between two frames:

$$\Delta_{gt}(\delta) = P(K + \delta) - P(K), \quad (2)$$

where  $P(t)$  is the recorded 6DoF pose of the object at frame  $t$ . If  $\|\Delta_{gt}(\delta)\| = 0$ , we assign frame  $K + \delta$  with the stationary label ( $c = 0$ ) to indicate that the human is yet to move the object; otherwise  $c = 1$ .

**Data split.** To test the model generalizability to unseen instances, we reserved all 471 samples of Chairs 3 and 4 (see Fig. 2) to form a *unseen instance test set*. We also randomly split the rest of the collected data (totalling 17,872 samples) into a training set of 12908 samples, a validation set of 3227 samples, and a *seen instance test set* of 1747 samples. Each object (except Chairs 3 and 4) and each action type appear in all three splits.

## V. EXPERIMENTS

Comparative evaluation is conducted to validate our network design and the use of object dynamic descriptors. We also examine the model generalizability to unseen objects from the six categories.

**Training.** We used the Adam optimizer [17] with a learning rate of 0.001 for training our network. A mini-batch of 32 samples was fed to the network during training. We trained the neural network models until they converged and showed the best performance on the validation set. All models were trained in a category-agnostic manner across all training data. Training and evaluation were conducted on a machine running Ubuntu 18.04 with a 1.80GHz Intel Xeon Silver 4108 CPU, and an NVIDIA RTX 2080Ti GPU.

**Evaluation metrics.** Given the ground-truth (GT) 6D pose changes of an object and the predictions produced by the methods, we measured the errors in pose translation ( $mm$ ) and pose rotation ( $10^{-3}rad$ ). We also adopted the mean per joint position error ( $mm$ ) with respect to the human skeleton (MPJPE-H) and the object keypoints (MPJPE-O).

**Ablation setting.** To justify our design choices, we also evaluate the performance of our network without using the proposed object dynamic descriptors (denoted *Ours w/o Desc*). Specifically, we remove the object dynamic descriptors from

the input to the object pre-processing branch and feed only the keypoint information to the weighted pooling block in Figure 3. The channel size of the FC layer predicting the object motion is adapted accordingly as well.

We further remove all the object information (both its dynamic descriptors and the geometry) and use solely the human skeleton motion in the first 10 frames to predict both the human and object motions in the last 10 frames, which we denoted as *BaseGCN*. In this setting, the input object representation, object encoder, and the weighted pooling block are deducted from our network.

**Comparison to SOTA methods.** We compare our model against two state-of-the-art (SOTA) methods (C-TE [3] and CAHMP [6]). All of the methods are given a 1-second observation (10 frames) of the human-object interaction, and required to predict the future 1-second (10 frames) joint motion of the human and the object. We are specifically interested in the prediction of object pose changes as accurate prediction can lead to desirable robot assistive operation in human-robot collaboration tasks.

### A. Quantitative analysis

Firstly, we evaluate our proposed method (HO-GCN), its variants (Ours w/o Desc and BaseGCN), and the comparing methods (C-TE and CAHMP) on the test set containing *seen objects*. Different methods are compared in terms of each future frame, short-term (0.1–0.5s) and long-term (0.6–1.0s) in Tab. II.

Comparisons show that our method consistently outperforms C-TE [3] which uses a convolution-based encoder-decoder structure for human-object joint motion prediction by a large margin. While the comparison with CAHMP shows it performs slightly better than our method in terms of human motion prediction, our method achieves a marginal improvement over CAHMP on the MPJPE-O metric in both short term and long term.

Secondly, we also test the different models on *two unseen objects*. In particular, Chair 3 (see Fig. 2) is an armchair, which is largely different from the other chairs in appearance and geometry. Note that the same keypoint configuration is used as shown in Fig. 2 for motion capture. Even without attaching markers to the armrests, we can achieve satisfactory results (see Fig. 5), showing that the abstracted geometry is sufficient for predicting the motion of a rigid body in our framework, justifying the use of 12 keypoints for all objects.

Tab. III shows the averaged performances over the 10 future frames. We can see that generalizing the trained models to unseen instances is challenging; compared to their performances in the all-seen setting, all models obtained worse results. Under this condition, our HO-GCN outperforms its variant (Ours w/o Desc and BaseGCN) by a relatively large margin on both chairs, showing the efficacy of the object dynamic descriptors for model generalization.

When comparing to CAHMP [6], we see that our method consistently achieves better results on different unseen chairs regarding the object motion prediction, showing the feasibility of applying our method to human-robot collaborative tasks,

Metrics	Methods	0.1s	0.2s	0.3s	0.4s	0.5s	0.6s	0.7s	0.8s	0.9s	1.0s	Short-term (0.1s-0.5s)	Long-term (0.6s-1.0s)	Mean
		Trans. Err.	C-TE	33.4	40.3	54.1	73.1	90.3	96.9	120.0	142.9	168.9	196.8	58.2
Trans. Err.	BaseGCN	34.7	38.6	48.3	65.1	67.1	92.4	106.9	145.6	164.6	184.4	50.8	138.8	94.8
Trans. Err.	Ours w/o Desc	37.4	43.2	53.8	68.0	79.0	94.0	105.9	129.2	137.6	145.9	56.3	122.5	89.4
Trans. Err.	Ours	22.4	29.0	39.4	47.2	62.1	72.5	80.6	136.3	148.8	163.1	<b>40.0</b>	<b>120.3</b>	<b>80.2</b>
Rot. Err.	C-TE	46.1	64.9	84.4	111.8	159.7	150.8	210.7	208.1	229.4	251.0	93.4	210.0	151.7
Rot. Err.	BaseGCN	46.5	54.8	64.2	85.5	115.6	130.4	159.9	190.9	198.4	227.9	73.3	181.5	127.4
Rot. Err.	Ours w/o Desc	52.2	63.1	94.0	108.2	122.6	140.2	154.5	178.3	190.8	202.2	88.0	173.2	130.6
Rot. Err.	Ours	33.7	46.4	60.7	82.8	86.0	100.5	115.7	191.0	204.9	223.8	<b>61.9</b>	<b>167.2</b>	<b>114.6</b>
MPJPE-O	C-TE	36.4	44.7	60.1	81.5	102.1	107.3	135.3	158.0	183.4	212.4	65.0	159.3	112.1
MPJPE-O	CAHMP	25.8	42.7	59.5	75.9	92.2	108.3	124.0	139.6	155.3	171.2	59.2	139.7	99.5
MPJPE-O	BaseGCN	37.5	42.0	53.0	71.2	76.6	101.7	117.8	157.9	175.5	198.3	56.1	150.2	103.2
MPJPE-O	Ours w/o Desc	42.0	47.1	60.1	74.5	86.1	101.8	114.2	138.0	147.1	156.1	62.0	131.5	96.7
MPJPE-O	Ours	24.0	31.6	42.5	52.5	67.2	78.3	88.0	147.2	160.5	175.9	<b>43.6</b>	<b>130.0</b>	<b>86.8</b>
MPJPE-H	C-TE	25.9	48.1	70.9	89.1	110.8	124.0	137.7	151.0	166.5	180.9	69.0	152.0	112.1
MPJPE-H	CAHMP	22.3	42.4	60.0	76.0	91.1	105.7	120.1	134.7	149.7	165.1	<b>58.3</b>	135.1	<b>96.7</b>
MPJPE-H	BaseGCN	27.8	50.8	73.4	92.0	108.1	123.2	138.1	143.0	159.0	166.2	70.4	145.9	108.2
MPJPE-H	Ours w/o Desc	24.0	45.9	65.7	82.6	97.3	110.2	122.5	134.1	144.4	154.5	63.1	<b>133.2</b>	98.2
MPJPE-H	Ours	23.9	45.8	65.8	84.6	100.6	115.4	129.4	140.0	151.9	163.6	64.1	140.1	102.1

TABLE II

COMPARISON OF PERFORMANCES OF DIFFERENT METHODS ON SEEN OBJECTS. THE UPPER BLOCK REPORTS THE PREDICTION ERRORS REGARDING THE 6DOF POSE (I.E., TRANSLATION AND ROTATION). THE BOTTOM BLOCK REPORTS THE MPJPE WITH RESPECT TO THE OBJECT AND THE HUMAN.

while the gap between Ours and CAHMP regarding the human motion prediction is narrowed. Both our method based on GCN and CAHMP based on RNN gain large improvement when compared to C-TE [3].

#### Will the object dynamic descriptors benefit the task?

Comparing our method to its variant (Ours w/o Desc) for ablation, we observe consistent performance gains regarding MPJPE-O, the translation and rotation errors in both short/long-term settings. From Tab. II, the gains for the *Seen Objects* test set are observed substantial in the short-term phase, where the motion state of the object is changing as it was just moved by the human actor. Consistent gains are observed in Tab. III for *Unseen Objects* test set too. Such reliable prediction is crucial for predicting the intents of how the human actor is going to manipulate the objects, and supports the use of our object dynamic descriptors for describing the intrinsic physical properties of an object.

**Can human skeleton data alone be used to predict the joint motion?** We also compare the variant of our method (Ours w/o Desc) to the baseline model (BaseGCN) that uses only the human skeleton data to predict the human-object joint motion. We see a small improvement obtained by our variant in the long-term. This may be attributed to that in the long-term the object motion is relatively large, and thus the benefit of using the object information (in this case the object keypoint positions) is more obvious.

#### B. Qualitative analysis

We visualized some randomly sampled results for qualitative analysis and comparison. A collection of qualitative results are shown in Fig. 4. Objects are displayed on the left. Results of the human and the object predicted by our HO-GCN are shown in purple and red, respectively. We visualize three frames of the input sequence and three prediction frames. In this setting, the prediction results are quite close to the GT. We show a failure case (observed across all methods) at the bottom right where the actor was lifting the small box (Object ID 7) and a large difference between GT and the prediction is seen. We assume that because the dataset mainly contains large-size objects such as chairs and tables, the trained model fails to generalize well to this smaller box, while producing a

Method	Tran. Err.	Rot. Err.	MPJPE-O	MPJPE-H	Object
C-TE	148.9	332.4	188.1	134.8	Chair 3
CAHMP	-	-	168.9	124.5	
BaseGCN	135.3	311.7	168.2	131.7	
Ours w/o Desc	129.6	304.1	160.8	<b>124.0</b>	
Ours	<b>122.1</b>	<b>296.5</b>	<b>151.0</b>	128.9	
C-TE	175.7	394.3	214.8	108.1	Chair 4
CAHMP	-	-	165.4	<b>85.8</b>	
BaseGCN	160.4	350.2	196.2	100.9	
Ours w/o Desc	131.0	319.2	164.1	91.7	
Ours	<b>124.2</b>	<b>310.8</b>	<b>156.7</b>	89.0	
C-TE	160.7	359.6	199.8	123.1	Sample Mean
CAHMP	-	-	167.4	<b>107.5</b>	
BaseGCN	146.4	328.6	180.5	118.2	
Ours w/o Desc	130.2	310.8	162.2	109.8	
Ours	<b>123.0</b>	<b>302.8</b>	<b>153.5</b>	111.4	

TABLE III

MODEL GENERALIZATION TEST ON UNSEEN OBJECTS *Chair 3* AND *Chair 4*. SPECIFICALLY, *Chair 3* IS AN ARMCHAIR WITH DRASTICALLY DIFFERENT GEOMETRY FROM CHAIRS (1 AND 2) FOR TRAINING.

satisfactory result for the medium-sized box (Object ID 8) for the *push forward* action.

Fig. 5 shows the qualitative comparison of different methods. Two input motion sequences, *Lift Box 8* and *Rotate Chair 4*, are shown at the top. The former sequence was sampled from the all-seen test set while the latter from the unseen. We show two predicted frames at  $t = 0.5s$  and  $1.0s$  corresponding to the short-term and long-term results. All methods produce reasonable results regarding the *Lift Box 8* sequence. However, C-TE and one of our variants (Our w/o Desc) seem to predict motions faster than the ground-truth. For *Rotate Chair 4*, only our method produces a plausible motion of rotating the chair. CAHMP, on the other hand, predicts a rough translation without rotating the chair, failing to capture the intention as the other methods do.

#### C. Human-robot collaborative tasks

We also showcase that the proposed HO-GCN can be useful for human-robot collaborative tasks. Specifically, we are interested in whether HO-GCN can successfully predict the manipulation intent of the human actor in terms of the future motion of the target object. We thus conducted experiments with a UR5 robot arm and a vacuum gripper with four silicone

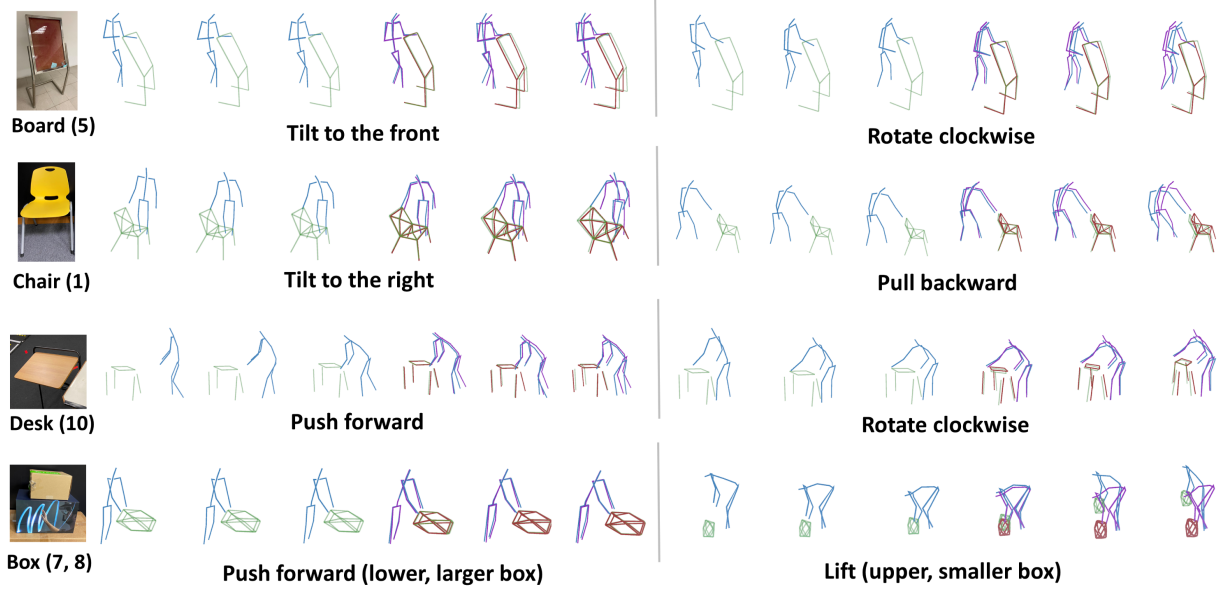


Fig. 4. Qualitative results produced by HO-GCN with the object dynamic descriptors are shown. Real-world objects are shown on the left. Results of the human and the object predicted by our HO-GCN are shown in purple and red, respectively, while GT is drawn in blue and green. Three frames of the input sequence and three prediction frames are visualized. A failure case is shown at the bottom right.

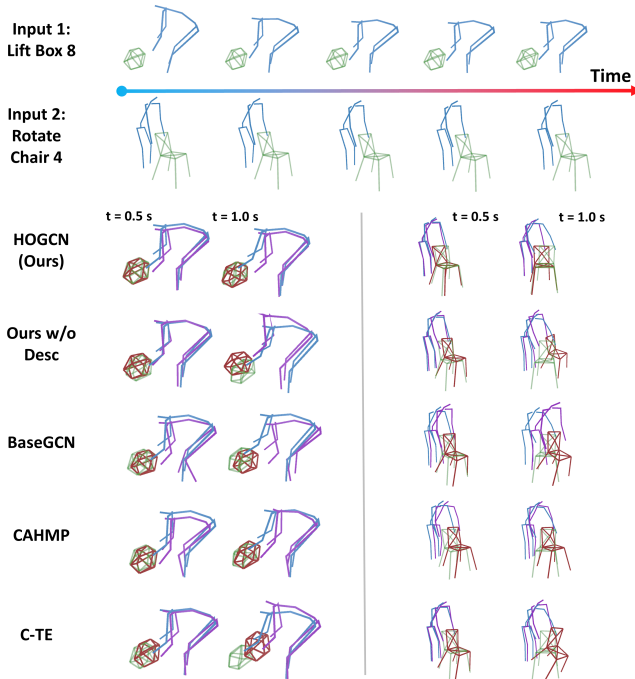


Fig. 5. Qualitative comparison of different methods. *Lift Box 8* is sampled from the seen object test set while *Rotate Chair 4* is from the unseen one. Results of the human and the object predicted by our HO-GCN are shown in purple and red, respectively, while GT's are drawn in blue and green. Our method can well capture the intended motion.

suction cups (of diameter 40mm). We chose a box which is suitable in weight and size for demonstration. The end-effector of the robot arm was pre-aligned to the target surface of the box. We used the OptiTrack system to provide 3D human motion information as well as the starting poses of the object.

With an input sequence of a human manipulating the box, our proposed method can predict the potential movement of the object in terms of its 6DoF changes. We then converted it to the corresponding robot trajectory that best produces the predicted object movement. In this showcase, we didn't consider the leading/following roles of the human and the robot in the cooperative task which could be our future work.

Our HO-GCN ran at a real-time rate (42 FPS) for processing such a 10-frame sequence on a desktop running Ubuntu 16.04 with an Intel Core i7-9700K CPU and an NVIDIA RTX 2080 GPU. Some results are shown in Fig. 6 and more can be found in our supplemental video.



Fig. 6. Experiment results of the collaborative tasks. The human actors were performing lifting, pushing, clockwise rotating, and counter-clockwise rotating from left to right.

## VI. CONCLUSIONS

In this paper, we focus on predicting full-body human interactions with large-sized daily objects and contribute to a large-scale dataset. Given as input a sequential observation of the human-object interaction, we design a novel graph convolutional network for predicting the future interactive motion. We also show that the object dynamic descriptors encoding the inherent physical properties of an object are beneficial to our network's generalization to unseen instances during test. Some showcasing examples of human-robot collaboration were presented using the proposed motion prediction method.

Currently, a naive version of the dynamic descriptor has been investigated; finding a more powerful representation for the inherent dynamic properties of an object is a promising future work. We would like to extend our human-robot collaboration to more objects and realistic settings. This would require extending our current work to predicting human-object interactions from video inputs.

## REFERENCES

- [1] “Optitrack motion tracking system [online],” in *Available: <http://www.naturalpoint.com/optitrack/>*.
- [2] P. Battaglia, R. Pascanu, M. Lai, D. J. Rezende, *et al.*, “Interaction networks for learning about objects, relations and physics,” in *Advances in neural information processing systems*, 2016, pp. 4502–4510.
- [3] J. Butepage, M. J. Black, D. Kragic, and H. Kjellstrom, “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.
- [4] Z. Cao, H. Gao, K. Mangalam, Q.-Z. Cai, M. Vo, and J. Malik, “Long-term human motion prediction with scene context,” *arXiv preprint arXiv:2007.03672*, 2020.
- [5] G. Cheng, K. Ramirez-Amaro, M. Beetz, and Y. Kuniyoshi, “Purposeful learning: Robot reasoning about the meanings of human activities,” *Science Robotics*, vol. 4, no. 26, 2019.
- [6] E. Corona, A. Pumarola, G. Alenya, and F. Moreno-Noguer, “Context-aware human motion prediction,” in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2020, pp. 6992–7001.
- [7] Q. Cui, H. Sun, and F. Yang, “Learning dynamic relationships for 3d human motion prediction,” in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2020, pp. 6519–6527.
- [8] Z. Ding, P. Wang, P. O. Ogunbona, and W. Li, “Investigation of different skeleton features for cnn-based 3d action recognition,” in *2017 IEEE International Conference on Multimedia & Expo Workshops (ICMEW)*. IEEE, 2017, pp. 617–622.
- [9] Y. Du, Y. Fu, and L. Wang, “Skeleton based action recognition with convolutional neural network,” in *2015 3rd IAPR Asian Conference on Pattern Recognition (ACPR)*. IEEE, 2015, pp. 579–583.
- [10] K. Ehsani, S. Tulsiani, S. Gupta, A. Farhadi, and A. Gupta, “Use the force, luke! learning to predict physical forces by simulating effects,” in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2020, pp. 224–233.
- [11] C. Finn, I. Goodfellow, and S. Levine, “Unsupervised learning for physical interaction through video prediction,” *arXiv preprint arXiv:1605.07157*, 2016.
- [12] G. Garcia-Hernando, S. Yuan, S. Baek, and T.-K. Kim, “First-person hand action benchmark with rgb-d videos and 3d hand pose annotations,” in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2018, pp. 409–419.
- [13] C. Ionescu, D. Papava, V. Olaru, and C. Sminchisescu, “Human3.6m: Large scale datasets and predictive methods for 3d human sensing in natural environments,” *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 36, no. 7, pp. 1325–1339, jul 2014.
- [14] M. Janner, S. Levine, W. T. Freeman, J. B. Tenenbaum, C. Finn, and J. Wu, “Reasoning about physical interactions with object-oriented prediction and planning,” *arXiv preprint arXiv:1812.10972*, 2018.
- [15] S. Kasahara, K. Konno, R. Owaki, T. Nishi, A. Takeshita, T. Ito, S. Kasuga, and J. Ushiba, “Malleable embodiment: Changing sense of embodiment by spatial-temporal deformation of virtual human body,” in *Proceedings of the 2017 CHI Conference on Human Factors in Computing Systems*, 2017, pp. 6438–6448.
- [16] S. Kim, K. Yun, J. Park, and J. Y. Choi, “Skeleton-based action recognition of people handling objects,” in *2019 IEEE Winter Conference on Applications of Computer Vision (WACV)*. IEEE, 2019, pp. 61–70.
- [17] D. P. Kingma and J. Ba, “Adam: A method for stochastic optimization,” *arXiv preprint arXiv:1412.6980*, 2014.
- [18] H. Koppula and A. Saxena, “Learning spatio-temporal structure from rgb-d videos for human activity detection and anticipation,” in *International conference on machine learning*, 2013, pp. 792–800.
- [19] H. S. Koppula and A. Saxena, “Anticipating human activities for reactive robotic response,” in *IROS*. Tokyo, 2013, p. 2071.
- [20] J.-F. Laffèche, S. Saunderson, and G. Nejat, “Robot cooperative behavior learning using single-shot learning from demonstration and parallel hidden markov models,” *IEEE Robotics and Automation Letters*, vol. 4, no. 2, pp. 193–200, 2018.
- [21] I. Lee, D. Kim, S. Kang, and S. Lee, “Ensemble deep learning for skeleton-based action recognition using temporal sliding lstm networks,” in *Proceedings of the IEEE International Conference on Computer Vision*, 2017, pp. 1012–1020.
- [22] C. Li, Q. Zhong, D. Xie, and S. Pu, “Co-occurrence feature learning from skeleton data for action recognition and detection with hierarchical aggregation,” *arXiv preprint arXiv:1804.06055*, 2018.
- [23] M. Li, S. Chen, X. Chen, Y. Zhang, Y. Wang, and Q. Tian, “Actional-structural graph convolutional networks for skeleton-based action recognition,” in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2019, pp. 3595–3603.
- [24] Z. Li, J. Sedlar, J. Carpentier, I. Laptev, N. Mansard, and J. Sivic, “Estimating 3d motion and forces of person-object interactions from monocular video,” in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2019, pp. 8640–8649.
- [25] J. Liu, A. Shahroudy, M. Perez, G. Wang, L.-Y. Duan, and A. C. Kot, “Ntu rgb+d 120: A large-scale benchmark for 3d human activity understanding,” *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2019.
- [26] C. Mandery, M. Plappert, J. Borras, and T. Asfour, “Dimensionality reduction for whole-body human motion recognition,” in *2016 19th International Conference on Information Fusion (FUSION)*. IEEE, 2016, pp. 355–362.
- [27] C. Mandery, O. Terlemez, M. Do, N. Vahrenkamp, and T. Asfour, “The kit whole-body human motion database,” in *International Conference on Advanced Robotics (ICAR)*, 2015, pp. 329–336.
- [28] F. Ofli, R. Chaudhry, G. Kurillo, R. Vidal, and R. Bajcsy, “Berkeley mhad: A comprehensive multimodal human action database,” in *2013 IEEE Workshop on Applications of Computer Vision (WACV)*. IEEE, 2013, pp. 53–60.
- [29] A. Shahroudy, J. Liu, T.-T. Ng, and G. Wang, “Ntu rgb+d: A large scale dataset for 3d human activity analysis,” in *IEEE Conference on Computer Vision and Pattern Recognition*, June 2016.
- [30] W. Sheng, A. Thobbi, and Y. Gu, “An integrated framework for human-robot collaborative manipulation,” *IEEE Transactions on Cybernetics*, vol. 45, no. 10, pp. 2030–2041, 2015.
- [31] C. Si, Y. Jing, W. Wang, L. Wang, and T. Tan, “Skeleton-based action recognition with spatial reasoning and temporal stack learning,” in *Proceedings of the European Conference on Computer Vision (ECCV)*, 2018, pp. 103–118.
- [32] O. Taheri, N. Ghorbani, M. J. Black, and D. Tzionas, “Grab: A dataset of whole-body human grasping of objects,” *arXiv preprint arXiv:2008.11200*, 2020.
- [33] Y. Tang, Y. Tian, J. Lu, P. Li, and J. Zhou, “Deep progressive reinforcement learning for skeleton-based action recognition,” in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2018, pp. 5323–5332.
- [34] D. Vogt, S. Stepputtis, S. Grehl, B. Jung, and H. Ben Amor, “A system for learning continuous human-robot interactions from human-human demonstrations,” in *2017 IEEE International Conference on Robotics and Automation (ICRA)*, 2017, pp. 2882–2889.
- [35] F. Xia, A. R. Zamir, Z. He, A. Sax, J. Malik, and S. Savarese, “Gibson env: Real-world perception for embodied agents,” in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2018, pp. 9068–9079.
- [36] L. Xia, C. Chen, and J. Aggarwal, “View invariant human action recognition using histograms of 3d joints,” in *Computer Vision and Pattern Recognition Workshops (CVPRW)*, 2012 IEEE Computer Society Conference on. IEEE, 2012, pp. 20–27.
- [37] S. Yan, Y. Xiong, and D. Lin, “Spatial temporal graph convolutional networks for skeleton-based action recognition,” in *Thirty-Second AAAI Conference on Artificial Intelligence*, 2018.
- [38] P. Zhang, C. Lan, J. Xing, W. Zeng, J. Xue, and N. Zheng, “View adaptive recurrent neural networks for high performance human action recognition from skeleton data,” in *Proceedings of the IEEE International Conference on Computer Vision*, 2017, pp. 2117–2126.
- [39] S. Zhang, Y. Zhang, Q. Ma, M. J. Black, and S. Tang, “Place: Proximity learning of articulation and contact in 3d environments.”
- [40] X. Zhao, S. Chumkamon, S. Duan, J. Rojas, and J. Pan, “Collaborative human-robot motion generation using lstm-rnn,” in *2018 IEEE-RAS 18th International Conference on Humanoid Robots (Humanoids)*, 2018, pp. 1–9.