

# Scaling Up Dynamic Human-Scene Interaction Modeling

Nan Jiang<sup>1,2\*</sup>, Zhiyuan Zhang<sup>1,2\*</sup>, Hongjie Li<sup>1</sup>, Xiaoxuan Ma<sup>3</sup>, Zan Wang<sup>4</sup>,  
Yixin Chen<sup>2</sup>, Tengyu Liu<sup>2</sup>, Yixin Zhu<sup>1</sup>✉, Siyuan Huang<sup>2</sup>✉

<sup>1</sup>Institute for AI, Peking University   <sup>2</sup>National Key Lab of General AI, BIGAI   <sup>3</sup>School of Computer Science, CPCS, Peking University  
<sup>4</sup>Beijing Institute of Technology   \*Equal contributors   ✉ yixin.zhu@pku.edu.cn, syhuang@bigai.ai

<https://jnnan.github.io/trumans/>

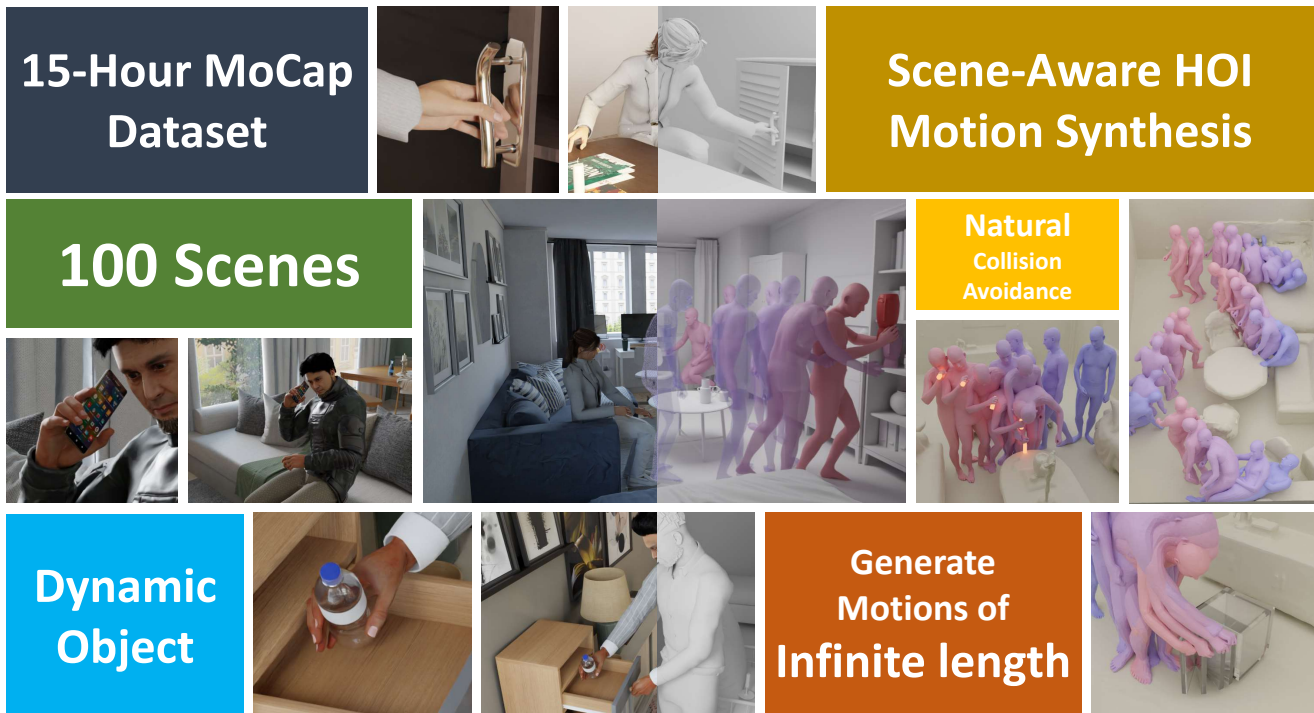


Figure 1. **Overview of TRUMANS dataset and our Human-Scene Interaction (HSI) framework.** We introduce the most extensive motion-captured HSI dataset, featuring diverse HSIs precisely captured in 100 scene configurations. Benefiting from TRUMANS, we propose a novel method for generating HSIs with arbitrary length, surpassing all baselines and exhibiting superb zero-shot generalizability.

## Abstract

Confronting the challenges of data scarcity and advanced motion synthesis in HSI modeling, we introduce the TRUMANS (*Tracking Human Actions in Scenes*) dataset alongside a novel HSI motion synthesis method. TRUMANS stands as the most comprehensive motion-captured HSI dataset currently available, encompassing over 15 hours of human interactions across 100 indoor scenes. It intricately captures whole-body human motions and part-level object dynamics, focusing on the realism of contact. This dataset is further scaled up by transforming physical environments into exact virtual models and applying extensive augmentations to appearance and motion for both humans and objects while maintaining interaction fidelity. Utilizing TRUMANS, we de-

vised a diffusion-based autoregressive model that efficiently generates Human-Scene Interaction (HSI) sequences of any length, taking into account both scene context and intended actions. In experiments, our approach shows remarkable zero-shot generalizability on a range of 3D scene datasets (e.g., PROX, Replica, ScanNet, ScanNet++), producing motions that closely mimic original motion-captured sequences, as confirmed by quantitative experiments and human studies.

## 1. Introduction

The intricate interplay between humans and their environment is a focal point in Human-Scene Interaction (HSI) [12], spanning diverse facets from object-level interaction [2, 25] to scene-level planning and interaction [1, 15, 16, 18]. While

significant strides have been made, the field is notably hindered by a scarcity of high-quality datasets. Early datasets like PiGraphs [39] and PROX [16] initiated the exploration but are constrained by scalability and data quality. MoCap datasets [14, 30] prioritize high-quality human motion capture using sophisticated equipment like VICON. However, they often lack in capturing diverse and immersive HSIs. Scalable datasets recorded via RGBD videos offer broader utility but are impeded by lower quality in human pose and object tracking. The advent of synthetic datasets [1, 3, 4, 56] provides cost efficiency and adaptability but fails to encapsulate the full spectrum of realistic HSIs, particularly in capturing dynamic 3D contacts and object tracking.

To address these challenges, this work first introduces the TRUMANS (Tracking Human Actions in Scenes) dataset. TRUMANS emerges as the most extensive motion-captured HSI dataset, **encompassing over 15 hours of diverse human interactions across 100 indoor scenes**. It captures whole-body human motions and part-level object dynamics with an emphasis on the realism of contact. This dataset is further enhanced by digitally replicating physical environments into accurate virtual models. Extensive augmentations in appearance and motion are applied to both humans and objects, ensuring high fidelity in interaction.

Next, we devise a computational model tackling the above challenges by taking both scene and action as conditions. Specifically, our model employs an autoregressive conditional diffusion with **scene** and **action** embeddings as conditional input, capable of generating motions of arbitrary length. To integrate **scene** context, we develop an efficient local scene perceiver by querying the global scene occupancy on a localized basis, which demonstrates robust proficiency in 3D-aware collision avoidance while navigating cluttered scenes. To incorporate frame-wise **action** labels as conditions, we integrate temporal features into action segments, empowering the model to accept instructions anytime while adhering to the given action labels. This dual integration of scene and action conditions enhances the controllability of our method, providing a nuanced interface for synthesizing plausible long-term motions in 3D scenes.

We conducted a comprehensive cross-evaluation of both the TRUMANS dataset and our motion synthesis method. Comparing TRUMANS with existing ones, we demonstrate that TRUMANS markedly improves the performance of current state-of-the-art approaches. Moreover, our method, evaluated both qualitatively and quantitatively, exceeds existing motion synthesis methods in terms of quality and zero-shot generalizability on unseen 3D scenes, closely approximating the quality of original motion-captured data. Beyond motion synthesis, TRUMANS has been benchmarked for human pose and contact estimation tasks, demonstrating its versatility and establishing it as a valuable asset for a broad range of future research endeavors.

Summarized in Fig. 1, our work significantly advances HSI modeling. Our contributions are threefold: (i) The introduction of TRUMANS, an extensive MoCap HSI dataset capturing a wide array of human behaviors across 100 indoor scenes, noted for its diversity, quality, and scalability. (ii) The development of a diffusion-based autoregressive method for the real-time generation of HSIs, adaptable to any length and conditioned on 3D scenes and action labels. (iii) Through extensive experimentation, we demonstrate the robustness of TRUMANS and our proposed methods, capable of generating motions that rival MoCap quality, outperforming existing baselines, and exhibiting exceptional zero-shot generalizability in novel environments.

## 2. Related Work

**HSI Datasets** Capturing human motions in 3D scenes is pivotal, with an emphasis on the quality and scale of human interactions. Early work focused on capturing coarse 3D human motions using 2D keypoints [33] or RGBD videos [39]. To improve quality and granularity, datasets like PROX [16] use scene scans as constraints to estimate SMPL-X parameters [36] from RGBD videos. However, these image-based motion capture methods often result in noisy 3D poses.

Recent efforts have incorporated more sophisticated systems like IMU or optical MoCap (*e.g.*, VICON) [14, 15, 17, 22, 30, 62], providing higher quality capture but limited in scalability. These are typically constrained to static scenes [15, 17, 56] or single objects [2, 22, 62], not fully representing complex real-world HSIs such as navigating cluttered spaces or managing concurrent actions.

Synthetic datasets [1, 4, 56] have attempted to fill this gap. Notable examples like BEDLAM [3] and CIRCLE [1] have been acknowledged for their cost efficiency and adaptability. These datasets integrate human motion data into synthetic scenes but fail to fully capture the range of realistic 3D HSIs, particularly in terms of dynamic object poses within their simulated environments.

Addressing these shortcomings, our work achieves a unique balance of quality and scalability. We replicate synthetic 3D environments in an optical motion capture setting, facilitating both accurate capture of humans and objects in complex HSIs and providing photorealistic renderings. This approach not only enhances the fidelity of the captured interactions but also extends the range of scenarios and environments that can be realistically simulated.

**HSI Generation** HSI generation involves single-frame human body [27, 61, 63] and temporal motion sequences [1, 17, 21, 26, 32, 35, 53–55, 58], utilizing models like conditional Variational Auto-Encoder (cVAE) [43] and diffusion models [19, 42, 44]. Recent advancements focus on generating arbitrary-length human motions through autoregressive methods [4, 7, 17, 31, 48, 60] and anchor frame genera-

Table 1. **Comparison of TRUMANS with existing HSI datasets.** TRUMANS differs by providing a diverse collection of HSIs, encompassing over 15 hours of interaction across 100 indoor scenes, along with photorealistic RGBD renderings in both multi-view and ego-view.

Datasets	Hours	MoCap	Human Representation	Dynamic Object	No. of Scenes	Contact Annotations	RGBD	Segmentation	Multi-view	Ego-view
GTA_IM [4]	9.3		skeleton		10		✓		✓	
PiGraphs [39]	2.0		skeleton		30		✓			
PROX [16]	0.9		SMPL-X		12	✓	✓	✓		
GRAB [47]	3.8	✓	SMPL-X	✓	-	✓				
SAMP [17]	1.7	✓	SMPL-X		-				✓	
RICH [20]	0.8		SMPL-X		5	✓	✓		✓	
BEHAVE [2]	4.2		SMPL	✓	-	✓	✓	✓	✓	
CHAIRS [22]	17.3	✓	SMPL-X	✓	-	✓	✓		✓	
COUCH [62]	3.0	✓	SMPL	✓	-	✓	✓	✓	✓	
iReplica [15]	0.8	✓	SMPL	✓	7	✓	✓		✓	✓
CIRCLE [1]	10.0	✓	SMPL-X		9					✓
TRUMANS	15.0	✓	SMPL-X	✓	100	✓	✓	✓	✓	✓

tion [37, 53]. Additionally, enhancing generation controllability has involved semantic guidance, such as action labels [64] and language descriptions [56, 57].

In comparison, our work contributes a conditional generative model with an autoregressive mechanism to generate **arbitrary-length** motions, combining diffusion model capabilities with improved **controllability** in HSI generation.

### 3. TRUMANS Dataset

This section introduces TRUMANS, the most comprehensive MoCap dataset dedicated to 3D HSIs thus far. TRUMANS offers not only accurate 3D ground truths but also photorealistic renderings accompanied by various 2D ground truths, suitable for various perceptual tasks in HSI. This section details the dataset’s statistics, data capture process, post-processing method, and our augmentation pipeline.

#### 3.1. Dataset Statistics

TRUMANS encompasses 15 hours of high-quality motion-captured data, featuring complex HSIs within 3D scenes, where humans interact with clustered environments and dynamic objects. Captured at a rate of 30 Hz using the state-of-the-art VICON MoCap system, the dataset comprises a total of 1.6 million frames. The HSI interactions in TRUMANS include 20 different types of common objects, ensuring a minimum of 5 distinct instances per type. The object categories encompass a range from static items like sofas and beds to dynamic objects such as bottles, and even articulated items including laptops and cabinets. TRUMANS incorporates performances from 7 participants (4 male and 3 female), who enacted various actions across 100 indoor scenes. These scenes span a variety of settings, such as dining rooms, living rooms, bedrooms, and kitchens, among others. For a comprehensive comparison of the TRUMANS dataset with existing HSI datasets, please refer to Tab. 1.

#### 3.2. Scene-aware Motion Capture

Aiming to capture realistic and diverse Human-Scene Interaction (HSI) within 3D scenes, our approach emphasizes both data quality and diversity. We initiate this process by replicating 3D scenes and objects sourced from the 3D-FRONT [10] dataset and BlenderKit [6] within the physical environment housing our MoCap devices. To ensure the naturalness of human interactions during motion capture, we meticulously create real-world placeholders that correspond to the affordances of the objects in the synthetic environment. All movable objects are tagged with markers compatible with the VICON system, enabling precise tracking of their poses. Actors undergo training to familiarize themselves with interacting with these placeholders. During the capturing sessions, actors are prompted to perform actions randomly selected from a pre-defined pool, ensuring a variety of interactions.

Post-capture, the human poses are converted into the SMPL-X format [36], employing a vertex-to-vertex optimization technique. This method is instrumental in calculating vertex-to-vertex distances between the human meshes and object meshes, facilitating accurate per-vertex contact annotations. We utilize Blender [5] to render multi-view photorealistic RGBD videos, segmentation masks, and ego-centric videos. To further diversify the renderings, we incorporate over 200 digital human models from Character Creator 4 [38], ensuring that objects strategically placed in scene backgrounds enhance the scene’s realism without impeding human movement. For a detailed exposition of our capture and processing pipeline, refer to Appendix B.4.

#### 3.3. MoCap Data Augmentation

Data augmentation is designed to adapt human motions to changes in 3D scene objects, ensuring physical plausibility and accuracy in HSI, following [45]. This process is vital in complex scenarios with concurrent or successive interac-

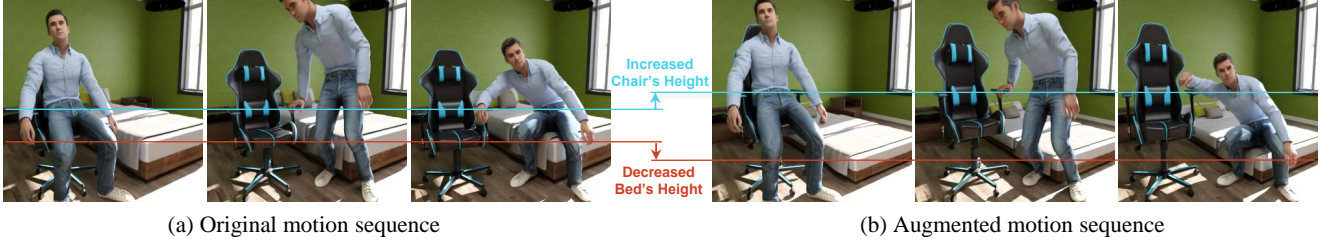


Figure 2. **Data augmentation for motion generation.** This example highlights how human motion is adjusted to accommodate variations in object sizes. Specifically, the chair’s height is increased, and the bed’s height is decreased, each by 15cm. Our augmentation method proficiently modifies human motion to maintain consistent interactions despite these changes in object dimensions.

tions; see Fig. 2. The pipeline consists of three main steps for integrating altered human motions into diverse 3D settings.

**Calculate Target Joint** We identify contact points between human joints and object meshes, and locate corresponding points on transformed or replaced objects. This step crucially adjusts the target joint’s position to maintain the original interaction’s contact relationship, ensuring realistic human-object interactions despite changes in object dimensions or positions.

**Refine Trajectory** To smooth out abrupt trajectory changes from the first step or Inverse Kinematic (IK) computations, we apply temporal smoothing to joint offsets, iteratively adjusting weights in adjacent frames. This refinement is critical for maintaining seamless motion, particularly in scenarios with multiple object interactions. Further details and theoretical background are discussed in Appendix B.5.

**Recompute Motion with IK** In the final step, we recompute human motion using the smoothed trajectories with an enhanced CCD-based [24] IK solver. This solver applies clipping and regularizations to bone movements, ensuring natural motion fluidity. Bones further from the root joint have increased rotational limits, reducing jitteriness and enhancing motion realism. For a complete description of these methods, refer to Appendix B.5.

## 4. Method

Utilizing the comprehensive TRUMANS dataset, we develop an autoregressive motion diffusion model. This model generates HSIs that are not only physically plausible in 3D scenes but also highly **controllable** through frame-wise action labels, capable of producing sequences of **arbitrary** length.

### 4.1. Problem Formulation and Notations

Given a 3D scene  $\mathcal{S}$ , a goal location  $\mathcal{G}$ , and action labels  $\mathcal{A}$ , our objective is to synthesize a human motion sequence  $\{\mathcal{H}_i\}_{i=1}^L$  of arbitrary length  $L$ . When interacting with dynamic objects  $\mathbf{P}$ , we also estimate the corresponding object pose sequence  $\{\mathcal{O}_i\}_{i=1}^L$ .

**Human** Human motion is represented as a sequence of parameterized human meshes  $\{\mathcal{H}_i\}$  using the SMPL-X model [36]. The motion is initially generated as body joints locations  $\{X^i\}_{i=1}^L$ , where  $X^i \in \mathbb{R}^{J \times 3}$  represents  $J = 24$  selected joints. These are fitted into the SMPL-X pose parameters  $\theta$ , global orientation  $\phi$ , hand poses  $h$ , and root translation  $r$ , resulting in the posed human mesh  $\mathcal{H} \in \mathbb{R}^{10475 \times 3}$ .

**Conditions** We formalize three types of conditions in our motion synthesis: 3D scene, goal location, and action labels. The 3D scene is represented by a voxel grid  $\mathcal{S} \in \{0, 1\}^{N_x \times N_y \times N_z}$ , with 1 indicating reachable locations. Goal locations are 2D positions  $\mathcal{G} \in \mathbb{R}^2$  for navigation, or 3D  $\mathbb{R}^3$  for joint-specific control. Action labels are multi-hot vectors  $\mathcal{A} \in \{0, 1\}^{L \times N_A}$ , indicating distinct actions.

**Object** When dynamic objects are involved, the object is represented by its point cloud  $\mathbf{P}$  in canonical coordinates and its global rotation  $R$  and translation  $T$ . The dynamic object sequence  $\{\mathcal{O}_i\}_{i=1}^L$  is then represented by sequences of rotations and translations  $\{R_i, T_i\}_{i=1}^L$ .

### 4.2. Autoregressive Motion Diffusion

Our model architecture is illustrated in Fig. 3. Our goal is to generate human motions that are not only physically plausible in 3D scenes but also highly controllable by frame-wise action labels, achieving arbitrary length in real time. We employ an autoregressive diffusion strategy where a long motion sequence is progressively generated by *episodes*, each defined as a motion segment of  $L_{epi}$  frames. Based on the approach by Shafir et al. [40], successive episodes are generated by extending from the final  $k$  frames of the prior episode. For each new episode, the first  $k$  frames are set based on the previous episode’s last  $k$  frames, with the noise on these transition frames zeroed out using a mask  $\mathbf{M}_{trans}$ . Our model aims to inpaint the remainder of each episode by filling in the unmasked frames.

To ensure precise control over character navigation and detailed interactions in each episode, we segment the overall goal  $\mathcal{G}$  into discrete subgoals, represented as  $\{\mathcal{G}_i\}_{i=1}^{N_{epi}}$ , where  $N_{epi}$  denotes the number of episodes. For navigation, each subgoal  $\mathcal{G}_i \in \mathbb{R}^2$  dictates the desired  $xy$ -coordinates of

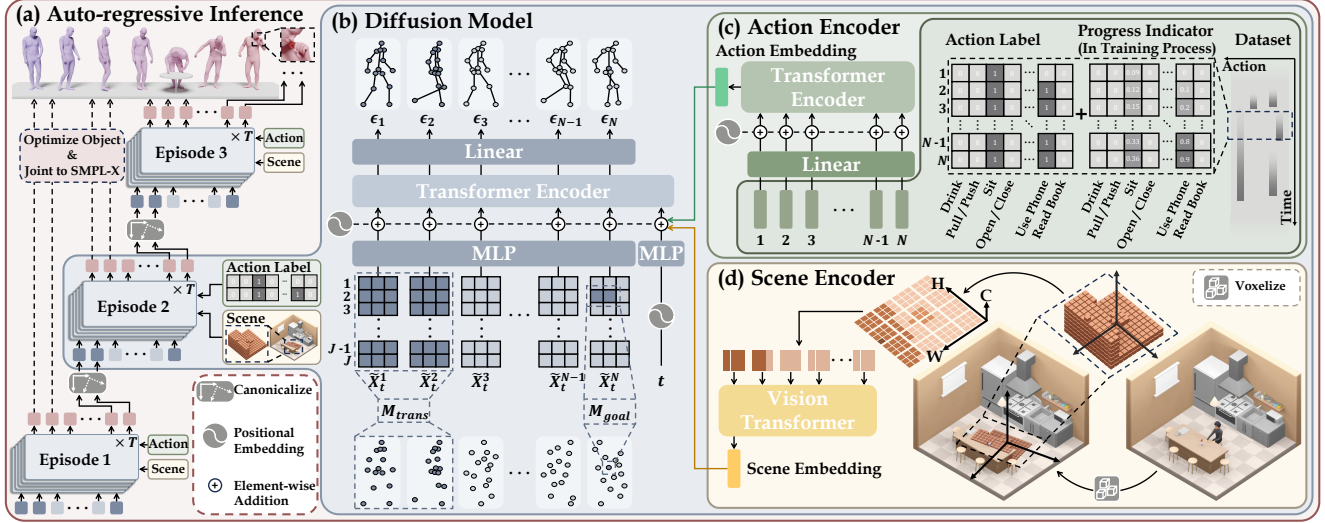


Figure 3. **Model architecture.** (a) Our model employs an autoregressive diffusion sampling approach to generate arbitrary long-sequence motions. (b) Within each episode, we synthesize motion using DDPM integrated with a transformer architecture, taking the human joint locations as input. (c)(d) Action and scene conditions are encoded and forwarded to the first token, guiding the motion synthesis process.

the character’s pelvis at an episode’s conclusion. Mirroring the masking approach used in  $\mathbf{M}_{trans}$ , we align the pelvis’s  $xy$ -coordinate in the episode’s final frame to the respective subgoal, simultaneously masking the corresponding diffusion noise. As the  $z$ -coordinate is unspecified, the model is trained to infer the appropriate pelvis height based on the scene setup, such as making the character sit when the subgoal indicates a chair’s location. This principle also governs fine-grained interactions, like grasping or pushing, where the subgoal  $\mathcal{G}_i \in \mathbb{R}^3$  is set to the precise 3D location, aligning the relevant hand joint to  $\mathcal{G}_i$  and masking joint noise accordingly. This specific masking on the subgoals is denoted as  $\mathbf{M}_{goal}$ . Motion

We devise a conditional diffusion model for generating motions within each episode. This process involves sampling from a Markov noising process  $\{X_t\}_{t=0}^T$ . Starting with the original human joint data  $X_0$  drawn from the data distribution, Gaussian noise is added to the components of  $X_0$  not masked by  $\mathbf{M} = \mathbf{M}_{trans} \cup \mathbf{M}_{goal}$ . The unmasked components, represented as  $(1 - \mathbf{M}) \odot X_t$  or  $\tilde{X}_t$  (where  $\odot$  is the Hadamard product), undergo a forward noising process

$$q(\tilde{X}_t | \tilde{X}_{t-1}) = \mathcal{N}(\tilde{X}_t; \sqrt{\alpha_t} \tilde{X}_{t-1}, (1 - \alpha_t)I), \quad (1)$$

with  $\alpha_t \in (0, 1)$  denoting hyper-parameters related to the variance schedule.

Motion data generation within our model employs a reversed diffusion process to gradually denoise  $\tilde{X}_T$ . Consistent with established diffusion model training methodologies, noise  $\epsilon_t$  is applied to obtain  $\tilde{X}_t$ , and a neural network  $\epsilon_\theta(\tilde{X}_t, t, \mathcal{S}, \mathcal{A})$  is constructed to approximate this noise. The

learning objective for  $\epsilon_\theta$  follows a simple objective [19]

$$\mathcal{L} = E_{\tilde{X}_0 \sim q(\tilde{X}_0 | \mathcal{C}), t \sim [1, T]} \left\| \epsilon - \epsilon_\theta(\tilde{X}_t, t, \mathcal{S}, \mathcal{A}) \right\|_2^2. \quad (2)$$

We adopt the Transformer model architecture [49], wherein the first token encodes information about the diffusion step, scene, and action, and subsequent tokens represent the noisy joint locations for each frame in the current episode. Throughout the sampling process, the model predicts the noise applied to each joint element. Once this sampling phase concludes, the joint locations are translated into SMPL-X parameters via a lightweight MLP. This translation is further refined through an optimization process, ensuring accurate alignment with the human joint data.

Upon generating the human motion sequence  $\{\mathcal{H}_i\}_{i=0}^L$ , we optimize the trajectory of the interacting object  $\{\mathcal{O}_i\}_{i=0}^L$  to ensure natural Human-Object Interactions (HOIs). To enhance the realism of the interaction, we further fine-tune the object’s pose in each frame to minimize the variance in distance between the object and the interacting hand [11].

### 4.3. Local Scene Perceiver

As illustrated in Fig. 3(d), the local scene perceiver is essential for embedding the local scene context, serving as a condition for motion generation. This component analyzes the scene using a local occupancy grid centered around the subgoal location for the current episode. Starting with the global occupancy grid  $\mathcal{S}$  of the scene, where each cell’s boolean value indicates reachability (1 for reachable, 0 otherwise), we focus on the  $i$ -th episode’s subgoal  $\mathcal{G}_i = (x, y, z)$  or  $(x, y)$ . A local occupancy grid is constructed around  $(x, y)$ , extending vertically from 0 to 1.8m. The grid’s orientation aligns with

the yaw of the agent’s pelvis at the episode’s start, and cell values are derived by querying the global occupancy grid.

The voxel grid is encoded using a Vision Transformer (ViT) [9]. We prepare the tokens by dividing the local occupancy grid into patches along the  $xy$ -plane, considering the  $z$ -axis as feature channels. These patches are then fed into the ViT model. The resulting scene embedding from this process is utilized as the condition for the diffusion model.

Discretizing the scene into a grid format is a necessary trade-off to boost training efficiency and practicality in our HSI method. Although directly generating the local occupancy grid from the scene mesh in real-time is technically feasible, it substantially prolongs training time. For instance, employing the *checksign* function from Kaolin results in a training process that is approximately 300 times slower, rendering it impractical. Despite this simplification, our empirical results demonstrate that the quality of motion generation is not significantly impacted by this approximation.

#### 4.4. Frame-wise Action Embedding

Our method distinguishes itself from prior approaches by incorporating frame-wise action labels into the long-term motion synthesis process, rather than generating a complete motion sequence from a singular action description. In our framework, a particular action can span multiple episodes, necessitating the model’s capability to comprehend the evolution and progression of an action over time.

To enhance our model’s understanding of action progression, we incorporate a progress indicator  $\mathcal{A}_{ind} \in \mathbb{R}^{L_{epi} \times N_A}$  into the frame-wise action labels, as depicted in Fig. 3(c). This indicator is realized by appending a real number  $n \in [0, 1]$  to the original action labels, representing the action’s advancement from start to finish. As a result, action labels take on values in  $0 \cup [1, 2]$  post-addition. For instance, during a drinking action from frame  $i$  to  $j$ , we modify the  $(0, 1)$  label by adding a value that linearly progresses from 0 to 1 across this interval. Thus, at the onset of drinking (frame  $i$ ), the label is augmented to 1, gradually increasing to 2 by frame  $j$ , the action’s conclusion. This nuanced labeling enables our model to seamlessly handle actions that span multiple episodes, significantly enhancing the realism and fluidity of the synthesized motion sequences.

The final action embedding is obtained by processing the progress-augmented action label  $\mathcal{A} \in \mathbb{R}^{L_{epi} \times N_A}$  through a Transformer encoder. Each frame’s action label  $\mathcal{A}_i \in \mathbb{R}^{N_A}$  is treated as an individual token in the Transformer’s input. The feature output from the last token is then passed through an MLP to generate the final action embedding.

## 5. Experiments

This section presents our evaluation of both TRUMANS and our proposed motion synthesis method, focusing on action-conditioned HSI generation. Additionally, we demonstrate

how TRUMANS contributes to advancements in state-of-the-art motion synthesis methods.

### 5.1. Experiment Settings

Our experimental evaluation of HSI generation quality is conducted under two distinct settings: *static* and *dynamic*. The *static* setting assesses synthesized motions in environments without dynamic interactable objects, concentrating on locomotion and interactions with static objects. Conversely, the *dynamic* setting evaluates motion synthesis involving interactions with dynamic objects. In both scenarios, we compare the performance of methods trained on TRUMANS with those trained on existing datasets [47, 63], offering a thorough insight into both the model’s efficacy and the dataset’s impact.

### 5.2. Baselines and Ablations

**Baselines–static setting** We compare TRUMANS with PROX [63], a dataset featuring human activities in indoor scenes. To ensure a fair comparison, we retain only the locomotion and scene interaction of static objects in TRUMANS, such as sitting and lying down. Baseline methods for this setting include cVAE [53], SceneDiff [21], and GMD [23].

**Baselines–dynamic setting** We compare TRUMANS with GRAB [47], known for capturing full-body grasping actions with human and object pose sequences. Here, the focus is on motions of interaction with dynamic objects, like drinking water and making phone calls, present in both datasets. We compare our method against IMoS [11] and GOAL [48], reproduced using their original implementations.

**Ablations** In our ablative studies, we examine the impact of disabling the action progress indicator  $\mathcal{A}_{ind}$  in our model. Additionally, to assess the significance of our data augmentation technique, we perform experiments using a non-augmented version of TRUMANS. For reference, our standard experiments employ the augmented TRUMANS, where each object is transformed into two different variations.

Our evaluation encompasses 10 unseen indoor scenes sourced from PROX [16], Replica [46], Scannet [8], and Scannet++ [59]. These scenes are adapted to the requirements of different methods, with modifications including conversion to point cloud format, voxelization, or maintaining their original mesh format. To evaluate the diversity of the synthesized motions, each method is tasked with generating five unique variations for each trajectory.

Furthermore, we conduct a qualitative comparison of our method with other recent approaches, such as SAMP [17], DIMOS [65], LAMA [25], and Wang et al. [55], based on the feasibility of reproducing these methods. Detailed findings from this comparison are discussed in Appendix A.4.

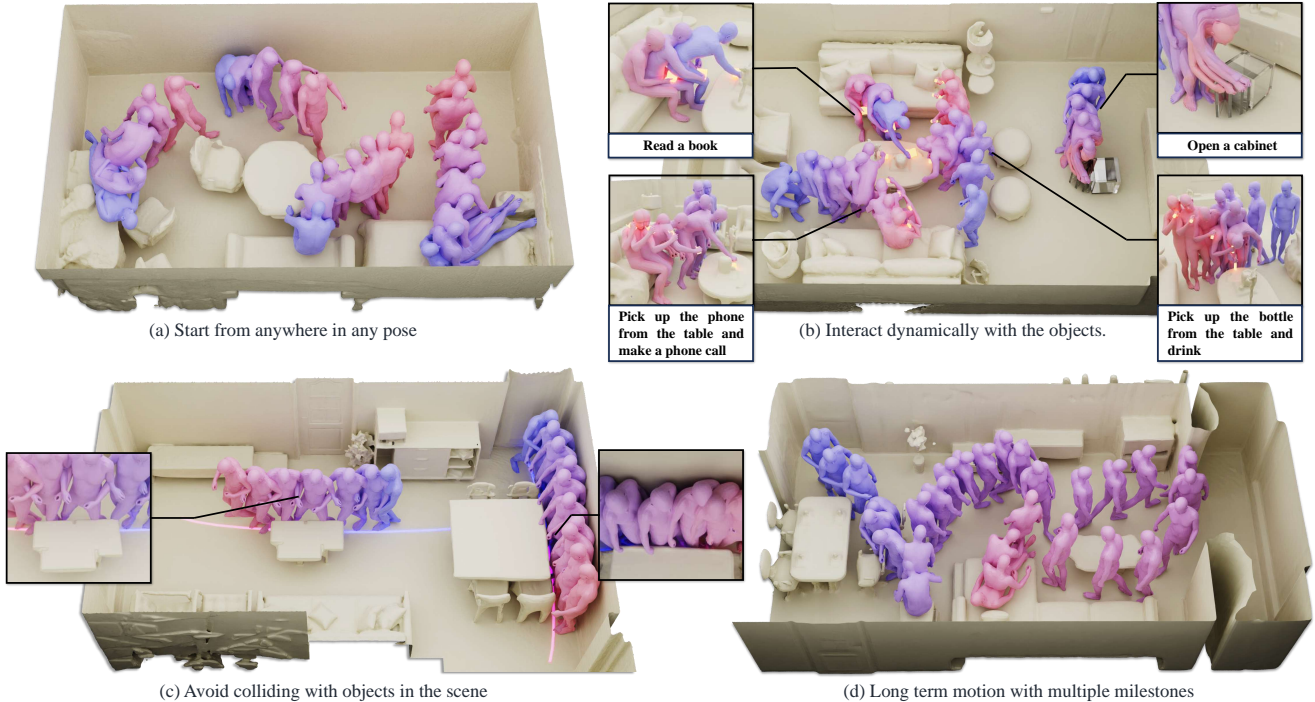


Figure 4. **Visualization of motion generation.** Leveraging local scene context and action instructions as conditions, our method demonstrates its proficiency in (a) initiating motion given the surrounding environment, (b) dynamically interacting with objects, (c) avoiding collisions during motion progression, and (d) robustly synthesizing long-term motion. The depicted scenes are selected from PROX, Replica, and FRONT3D-test datasets, none of which were included in the training phase. For qualitative results, please refer to the *Supplementary Video*.

### 5.3. Evaluation Metrics

In the *static* setting, we employ *Contact* and *Penetration* metrics, as recommended by Zhao et al. [65], to evaluate foot slide and object penetration issues in synthesized motions. These metrics measure the degree to which the synthesized motions conform to the specified scene. For the *dynamic* setting, we utilize *FID* and *Diversity* metrics, commonly used in language and action-guided motion generation tasks [11, 49]. These metrics measure the quality and diversity of HOI motion generation involving various small objects.

Additionally, we introduce a novel MoCap-differentiating human study for evaluation. Participants are presented with five sequences, one of which is motion-captured, and are asked to identify the MoCap sequence. The likelihood of correctly identifying the MoCap sequence serves as an indicator of the synthesized motion’s realism. We quantify this aspect through the Success Rate of Discrimination (SucRate-Dis), reflecting the percentage of participants who accurately identify the MoCap sequence.

### 5.4. Results and Analysis

Fig. 4 showcases our method’s qualitative strengths. It adeptly manages complex scene configurations, including initiating context-aware motion, avoiding collisions during movement, and generating extended motions, especially in

HOI scenarios involving dynamic object interaction.

In the *static* setting (Tab. 2), our method, trained on TRUMANS, surpasses baselines across most metrics. Notably, disabling data augmentation leads to increased penetration, suggesting the efficacy of augmented data in producing physically plausible motions. Compared to models trained on PROX, ours shows significant improvements, highlighting TRUMANS as a high-quality resource for HSI research.

Table 2. **Evaluation of locomotion and scene-level interaction.** We compare performances on TRUMANS and PROX [16].

Method	Cont.↑	Pene <sub>mean</sub> ↓	Pene <sub>max</sub> ↓	Dis. suc.↓
Wang et al. [53]	0.969	1.935	14.33	0.581
SceneDiff [21]	0.912	<b>1.691</b>	17.48	0.645
GMD [23]	0.931	2.867	21.30	0.871
Ours	<b>0.992</b>	1.820	<b>11.74</b>	<b>0.258</b>
Ours w/o aug.	0.991	2.010	15.52	-
Wang et al. [53]	0.688	4.935	34.10	0.903
SceneDiff [21]	0.712	3.267	27.48	0.935
GMD [23]	0.702	4.867	38.30	0.968
Ours	0.723	4.820	31.74	0.903

Tab. 3 illustrates results in the *dynamic* setting, where our approach excels in 3D HOI generation. High penetration rates with GRAB-trained methods indicate its limitations

in scene-adherent HOI motions, while TRUMANS captures more detailed interactions. The absence of the progress indicator  $\mathcal{A}_{ind}$  leads to method failure, as evidenced by the ablation study.

Table 3. **Evaluation of object-level interaction.** We compare performances on TRUMANS and GRAB [47]. The definition of “Real” follows the one defined in Tevet et al. [49]

Method	FID↓	Diversity→	Pene <sub>scene</sub> ↓	Dis. suc.↓
Real-TRUMANS	-	2.734	-	-
GOAL [48]	0.512	2.493	34.10	0.801
IMoS [11]	0.711	2.667	37.48	0.774
Ours	<b>0.313</b>	<b>2.693</b>	11.74	<b>0.226</b>
Ours - $\mathcal{A}_{ind}$	2.104	1.318	<b>10.62</b>	1.000
Real-GRAB [47]	-	2.155	-	-
GOAL [48]	0.429	2.180	44.09	0.801
IMoS [11]	0.410	2.114	41.50	0.774
Ours	0.362	2.150	34.41	0.516

Human studies further affirm the quality of our method. Only about a quarter of participants could distinguish our synthesized motions from real MoCap data, nearly aligning with the 1/5 SucRateDis of random guessing. This suggests that our synthesized motions are nearly indistinguishable from high-quality MoCap data. Comparative evaluations with recent methods [17, 25, 55, 65] show our model’s superiority, outperforming the second-best model by over 30% in support rate. For more detailed results, please refer to the *Supplementary Video*.

**Real-time Control** Our method can sample an episode of motion (1.6 seconds at 10 FPS) in 0.7 seconds on an A800 GPU. This efficiency enables uninterrupted long-term motion generation with a consistent control signal. For new control signals, to minimize the 0.7-second delay, we implement an incremental sampling strategy: initially, 2 frames are sampled immediately, followed by sampling 4 frames during their execution, increasing exponentially until 16 frames are sampled. This approach ensures a balance between real-time control and smooth motion continuity. Please refer to our *Supplementary Video* for a visual demonstration.

### 5.5. Additional Image-based Tasks

TRUMANS, with its photo-realistic renderings and per-vertex 3D contact annotations, is also suited for various image-based tasks. We focus on its application in 3D human mesh estimation and contact estimation.

**3D Human Mesh Estimation** For reconstructing 3D human body meshes from input images, we utilize the state-of-the-art method [29] as a baseline. We evaluate if including TRUMANS in training enhances performance on the 3DPW dataset [51]. Following Ma et al. [29], we report MPJPE, PA-MPJPE, and MPVE for the estimated poses and meshes.

**3D Contact Estimation** This task involves predicting per-vertex 3D contact on the SMPL mesh [28] from an input image. We compare TRUMANS against RICH [20] and DAMON [50], both featuring vertex-level 3D contact labels with RGB images. Utilizing BSTRO [20] for RICH and DECO [50] for DAMON, we measure precision, recall, F1 score, and geodesic error following the literature [20, 50].

**Results and Analysis** Quantitative results in Tab. 4 reveal that integrating TRUMANS with 3DPW significantly improves human mesh estimation. Contact estimation outcomes, presented in Tab. 5, show enhanced performance with TRUMANS, particularly in reducing geodesic error. These results suggest that combining synthetic data from TRUMANS with real-world data substantially benefits image-based tasks. For detailed experimental insights, see Appendix A.5.

Table 4. **Performance of Ma et al. [29] trained on 3DPW [51] combined with TRUMANS in different ratios.**

Training Data	MPVE↓	MPJPE↓	PA-MPJPE↓
3DPW [51]	101.3	88.2	54.4
3DPW+T (2:1)	88.8	<b>77.2</b>	<b>46.4</b>
3DPW+T (1:1)	<b>78.5</b>	78.5	<b>46.4</b>

Table 5. **Performance of BSTRO [20] and DECO [50] trained on RICH [20] and DAMON [50] combined with TRUMANS, respectively.**

Training Data	Prec↑	Rec↑	F1↑	geo err↓
RICH [20]	0.6823	<b>0.7427</b>	0.6823	10.27
R+T (2:1)	0.7087	0.7370	<b>0.6927</b>	9.593
R+T (1:1)	<b>0.7137</b>	0.7286	0.6923	<b>9.459</b>
DAMON [50]	0.6388	0.5232	0.5115	25.06
D+T (2:1)	0.6472	<b>0.5237</b>	<b>0.5148</b>	21.54
D+T (1:1)	<b>0.6701</b>	0.4806	0.4972	<b>18.87</b>

## 6. Conclusion

We introduce TRUMANS, a large-scale mocap dataset, alongside a novel motion synthesis method, addressing scalability, data quality, and advanced motion synthesis challenges in HSI modeling. As the most comprehensive dataset in its category, TRUMANS encompasses diverse human interactions with dynamic and articulated objects within 100 indoor scenes. Our diffusion-based autoregressive motion synthesis method, leveraging TRUMANS, is capable of real-time generation of HSI sequences of arbitrary length. Experimental results indicate that the motions generated by our method closely mirror the quality of the original MoCap data.

**Acknowledgment** The authors would like to thank NVIDIA for their generous support of GPUs and hardware. This work is supported in part by the National Science and Technology Major Project (2022ZD0114900) and the Beijing Nova Program.



## References

- [1] Joao Pedro Araújo, Jiaman Li, Karthik Vetrivel, Rishi Agarwal, Jiajun Wu, Deepak Gopinath, Alexander William Clegg, and Karen Liu. Circle: Capture in rich contextual environments. In *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2023. **1, 2, 3**
- [2] Bharat Lal Bhatnagar, Xianghui Xie, Ilya A Petrov, Cristian Sminchisescu, Christian Theobalt, and Gerard Pons-Moll. Behave: Dataset and method for tracking human object interactions. In *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2022. **1, 2, 3, A2**
- [3] Michael J Black, Priyanka Patel, Joachim Tesch, and Jinlong Yang. Bedlam: A synthetic dataset of bodies exhibiting detailed lifelike animated motion. In *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2023. **2**
- [4] Zhe Cao, Hang Gao, Karttikeya Mangalam, Qi-Zhi Cai, Minh Vo, and Jitendra Malik. Long-term human motion prediction with scene context. In *European Conference on Computer Vision (ECCV)*, 2020. **2, 3**
- [5] Blender Online Community. Blender - a 3d modelling and rendering package, 2018. **3**
- [6] Blender Online Community. Blenderkit. <https://www.blenderkit.com/>, 2023. **3, A3**
- [7] Enric Corona, Albert Pumarola, Guillem Alenyà, and Francesc Moreno-Noguer. Context-aware human motion prediction. In *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2020. **2**
- [8] Angela Dai, Angel X Chang, Manolis Savva, Maciej Halber, Thomas Funkhouser, and Matthias Nießner. Scannet: Richly-annotated 3d reconstructions of indoor scenes. In *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2017. **6**
- [9] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, Jakob Uszkoreit, and Neil Houlsby. An image is worth 16x16 words: Transformers for image recognition at scale. In *International Conference on Learning Representations (ICLR)*, 2021. **6**
- [10] Huan Fu, Bowen Cai, Lin Gao, Ling-Xiao Zhang, Jiaming Wang, Cao Li, Qixun Zeng, Chengyue Sun, Rongfei Jia, Binqiang Zhao, et al. 3d-front: 3d furnished rooms with layouts and semantics. In *International Conference on Computer Vision (ICCV)*, 2021. **3, A3**
- [11] Anindita Ghosh, Rishabh Dabral, Vladislav Golyanik, Christian Theobalt, and Philipp Slusallek. Imos: Intent-driven full-body motion synthesis for human-object interactions. In *Computer Graphics Forum*, 2023. **5, 6, 7, 8**
- [12] James J Gibson. *The perception of the visual world*. Houghton Mifflin, 1950. **1**
- [13] Chuan Guo, Shihao Zou, Xinxin Zuo, Sen Wang, Wei Ji, Xingyu Li, and Li Cheng. Generating diverse and natural 3d human motions from text. In *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2022. **A1**
- [14] Vladimir Guzov, Aymen Mir, Torsten Sattler, and Gerard Pons-Moll. Human poseitoning system (hps): 3d human pose estimation and self-localization in large scenes from body-mounted sensors. In *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2021. **2**
- [15] Vladimir Guzov, Julian Chibane, Riccardo Marin, Yannan He, Yunus Saracoglu, Torsten Sattler, and Gerard Pons-Moll. Interaction replica: Tracking human-object interaction and scene changes from human motion. In *International Conference on 3D Vision (3DV)*, 2023. **1, 2, 3, A2**
- [16] Mohamed Hassan, Vasileios Choutas, Dimitrios Tzionas, and Michael J Black. Resolving 3d human pose ambiguities with 3d scene constraints. In *International Conference on Computer Vision (ICCV)*, 2019. **1, 2, 3, 6, 7**
- [17] Mohamed Hassan, Duygu Ceylan, Ruben Villegas, Jun Saito, Jimei Yang, Yi Zhou, and Michael Black. Stochastic scene-aware motion prediction. In *International Conference on Computer Vision (ICCV)*, 2021. **2, 3, 6, 8, A1, A2**
- [18] Mohamed Hassan, Partha Ghosh, Joachim Tesch, Dimitrios Tzionas, and Michael J Black. Populating 3d scenes by learning human-scene interaction. In *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2021. **1**
- [19] Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2020. **2, 5**
- [20] Chun-Hao P Huang, Hongwei Yi, Markus Höschle, Matvey Safroshkin, Tsvetelina Alexiadis, Senya Polikovsky, Daniel Scharstein, and Michael J Black. Capturing and inferring dense full-body human-scene contact. In *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2022. **3, 8, A2**
- [21] Siyuan Huang, Zan Wang, Puhao Li, Baoxiong Jia, Tengyu Liu, Yixin Zhu, Wei Liang, and Song-Chun Zhu. Diffusion-based generation, optimization, and planning in 3d scenes. In *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2023. **2, 6, 7, A1**
- [22] Nan Jiang, Tengyu Liu, Zhexuan Cao, Jieming Cui, Zhiyuan Zhang, Yixin Chen, He Wang, Yixin Zhu, and Siyuan Huang. Full-body articulated human-object interaction. In *International Conference on Computer Vision (ICCV)*, 2023. **2, 3**
- [23] Korrawe Karunratanakul, Konpat Preechakul, Supasorn Suwajanakorn, and Siyu Tang. Guided motion diffusion for controllable human motion synthesis. In *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2023. **6, 7, A1**
- [24] Ben Kenwright. Inverse kinematics-cyclic coordinate descent (ccd). *Journal of Graphics Tools*, 2012. **4**
- [25] Jiye Lee and Hanbyul Joo. Locomotion-action-manipulation: Synthesizing human-scene interactions in complex 3d environments. In *International Conference on Computer Vision (ICCV)*, 2023. **1, 6, 8, A1**
- [26] Jiaman Li, Jiajun Wu, and C Karen Liu. Object motion guided human motion synthesis. *arXiv preprint arXiv:2309.16237*, 2023. **2**
- [27] Xueting Li, Sifei Liu, Kihwan Kim, Xiaolong Wang, Ming-Hsuan Yang, and Jan Kautz. Putting humans in a scene: Learning affordance in 3d indoor environments. In *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2019. **2**
- [28] Matthew Loper, Naureen Mahmood, Javier Romero, Gerard Pons-Moll, and Michael J. Black. SMPL: A skinned multi-

- person linear model. *ACM Transactions on Graphics (TOG)*, 2015. 8
- [29] Xiaoxuan Ma, Jiajun Su, Chunyu Wang, Wentao Zhu, and Yizhou Wang. 3d human mesh estimation from virtual markers. In *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2023. 8
- [30] Christian Mandery, Ömer Terlemez, Martin Do, Nikolaus Vahrenkamp, and Tamim Asfour. The kit whole-body human motion database. In *International Conference on Robotics and Automation (ICRA)*, 2015. 2
- [31] Wei Mao, Miaomiao Liu, Richard Hartley, and Mathieu Salzmann. Contact-aware human motion forecasting. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2022. 2
- [32] Aymen Mir, Xavier Puig, Angjoo Kanazawa, and Gerard Pons-Moll. Generating continual human motion in diverse 3d scenes. In *International Conference on 3D Vision (3DV)*, 2023. 2
- [33] Aron Monszpart, Paul Guerrero, Duygu Ceylan, Ersin Yumer, and Niloy J Mitra. imapper: interaction-guided scene mapping from monocular videos. *ACM Transactions on Graphics (TOG)*, 2019. 2
- [34] Gyeongsik Moon and Kyoung Mu Lee. I2l-meshnet: Image-to-lixel prediction network for accurate 3d human pose and mesh estimation from a single rgb image. In *European Conference on Computer Vision (ECCV)*, 2020. A2
- [35] Liang Pan, Jingbo Wang, Buzhen Huang, Junyu Zhang, Haofan Wang, Xu Tang, and Yangang Wang. Synthesizing physically plausible human motions in 3d scenes. In *International Conference on 3D Vision (3DV)*, 2023. 2
- [36] Georgios Pavlakos, Vasileios Choutas, Nima Ghorbani, Timo Bolkart, Ahmed AA Osman, Dimitrios Tzionas, and Michael J Black. Expressive body capture: 3d hands, face, and body from a single image. In *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2019. 2, 3, 4, A4
- [37] Huaijin Pi, Sida Peng, Minghui Yang, Xiaowei Zhou, and Hujun Bao. Hierarchical generation of human-object interactions with diffusion probabilistic models. In *International Conference on Computer Vision (ICCV)*, 2023. 3
- [38] Reallusion. Character creator 4. <https://www.reallusion.com/character-creator/>, 2023. 3, A3
- [39] Manolis Savva, Angel X Chang, Pat Hanrahan, Matthew Fisher, and Matthias Nießner. Pigraphs: learning interaction snapshots from observations. *ACM Transactions on Graphics (TOG)*, 2016. 2, 3
- [40] Yonatan Shafir, Guy Tevet, Roy Kapon, and Amit H Bermano. Human motion diffusion as a generative prior. *arXiv preprint arXiv:2303.01418*, 2023. 4
- [41] Vicon Software. Shogun. <https://www.vicon.com/software/shogun/>, 2023. A3
- [42] Jascha Sohl-Dickstein, Eric Weiss, Niru Maheswaranathan, and Surya Ganguli. Deep unsupervised learning using nonequilibrium thermodynamics. In *International Conference on Machine Learning (ICML)*, 2015. 2
- [43] Kihyuk Sohn, Honglak Lee, and Xinchen Yan. Learning structured output representation using deep conditional generative models. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2015. 2
- [44] Yang Song and Stefano Ermon. Generative modeling by estimating gradients of the data distribution. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2019. 2
- [45] Sebastian Starke, He Zhang, Taku Komura, and Jun Saito. Neural state machine for character-scene interactions. *ACM Transactions on Graphics (TOG)*, 38(6):209–1, 2019. 3
- [46] Julian Straub, Thomas Whelan, Lingni Ma, Yufan Chen, Erik Wijmans, Simon Green, Jakob J Engel, Raul Mur-Artal, Carl Ren, Shobhit Verma, et al. The replica dataset: A digital replica of indoor spaces. *arXiv preprint arXiv:1906.05797*, 2019. 6
- [47] Omid Taheri, Nima Ghorbani, Michael J Black, and Dimitrios Tzionas. Grab: A dataset of whole-body human grasping of objects. In *European Conference on Computer Vision (ECCV)*, 2020. 3, 6, 8
- [48] Omid Taheri, Vasileios Choutas, Michael J. Black, and Dimitrios Tzionas. GOAL: Generating 4D whole-body motion for hand-object grasping. In *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2022. 2, 6, 8, A1
- [49] Guy Tevet, Sigal Raab, Brian Gordon, Yoni Shafir, Daniel Cohen-or, and Amit Haim Bermano. Human motion diffusion model. In *International Conference on Learning Representations (ICLR)*, 2022. 5, 7, 8
- [50] Shashank Tripathi, Agniv Chatterjee, Jean-Claude Passy, Hongwei Yi, Dimitrios Tzionas, and Michael J Black. Deco: Dense estimation of 3d human-scene contact in the wild. In *International Conference on Computer Vision (ICCV)*, 2023. 8, A2
- [51] Timo Von Marcard, Roberto Henschel, Michael J Black, Bodo Rosenhahn, and Gerard Pons-Moll. Recovering accurate 3d human pose in the wild using imus and a moving camera. In *European Conference on Computer Vision (ECCV)*, 2018. 8, A2
- [52] Dongkai Wang and Shiliang Zhang. 3d human mesh recovery with sequentially global rotation estimation. In *International Conference on Computer Vision (ICCV)*, 2023. A2
- [53] Jiashun Wang, Huazhe Xu, Jingwei Xu, Sifei Liu, and Xiaolong Wang. Synthesizing long-term 3d human motion and interaction in 3d scenes. In *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2021. 2, 3, 6, 7, A1
- [54] Jingbo Wang, Sijie Yan, Bo Dai, and Dahua Lin. Scene-aware generative network for human motion synthesis. In *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2021.
- [55] Jingbo Wang, Yu Rong, Jingyuan Liu, Sijie Yan, Dahua Lin, and Bo Dai. Towards diverse and natural scene-aware 3d human motion synthesis. In *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2022. 2, 6, 8, A1
- [56] Zan Wang, Yixin Chen, Tengyu Liu, Yixin Zhu, Wei Liang, and Siyuan Huang. Humanise: Language-conditioned human motion generation in 3d scenes. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2022. 2, 3
- [57] Zeqi Xiao, Tai Wang, Jingbo Wang, Jinkun Cao, Wenwei Zhang, Bo Dai, Dahua Lin, and Jiangmiao Pang. Unified human-scene interaction via prompted chain-of-contacts. *arXiv preprint arXiv:2309.07918*, 2023. 3

- [58] Sirui Xu, Zhengyuan Li, Yu-Xiong Wang, and Liang-Yan Gui. InterDiff: Generating 3d human-object interactions with physics-informed diffusion. In *International Conference on Computer Vision (ICCV)*, 2023. [2](#)
- [59] Chandan Yeshwanth, Yueh-Cheng Liu, Matthias Nießner, and Angela Dai. Scannet++: A high-fidelity dataset of 3d indoor scenes. In *International Conference on Computer Vision (ICCV)*, 2023. [6](#)
- [60] Jianrong Zhang, Yangsong Zhang, Xiaodong Cun, Shaoli Huang, Yong Zhang, Hongwei Zhao, Hongtao Lu, and Xi Shen. T2m-gpt: Generating human motion from textual descriptions with discrete representations. In *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2023. [2](#)
- [61] Siwei Zhang, Yan Zhang, Qianli Ma, Michael J Black, and Siyu Tang. Generating person-scene interactions in 3d scenes. In *International Conference on 3D Vision (3DV)*, 2020. [2](#)
- [62] Xiaohan Zhang, Bharat Lal Bhatnagar, Sebastian Starke, Vladimir Guzov, and Gerard Pons-Moll. Couch: Towards controllable human-chair interactions. In *European Conference on Computer Vision (ECCV)*, 2022. [2](#), [3](#)
- [63] Yan Zhang, Mohamed Hassan, Heiko Neumann, Michael J Black, and Siyu Tang. Generating 3d people in scenes without people. In *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2020. [2](#), [6](#)
- [64] Kaifeng Zhao, Shaofei Wang, Yan Zhang, Thabo Beeler, and Siyu Tang. Compositional human-scene interaction synthesis with semantic control. In *European Conference on Computer Vision (ECCV)*, 2022. [3](#)
- [65] Kaifeng Zhao, Yan Zhang, Shaofei Wang, Thabo Beeler, and Siyu Tang. Synthesizing diverse human motions in 3d indoor scenes. In *International Conference on Computer Vision (ICCV)*, 2023. [6](#), [7](#), [8](#), [A1](#)
- [66] Yi Zhou, Connelly Barnes, Jingwan Lu, Jimei Yang, and Hao Li. On the continuity of rotation representations in neural networks. In *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2019. [A1](#)

## A. Additional Details of Experiments

This section offers a detailed overview of our experimental setup, including the implementation specifics of our method and the necessary adaptations made to baseline methods to ensure fair comparisons. We also elaborate on how our sampling strategy enables real-time control over character motion. For additional qualitative insights and extensive zero-shot transfer experiments, we direct readers to the accompanying *Supplementary Video*.

### A.1. Experiment Settings

Our experimental setup for motion synthesis methods includes two distinct settings: *static* and *dynamic*. In the *static* setting, we focus on evaluating the quality of locomotion and scene-level interactions. Each test scene features five predefined pairs of start and end points, given as  $(x, y)$  coordinates in a z-up world coordinate system. These points, often located on furniture like chairs, test the method’s ability to produce scene-appropriate motions. For complete trajectories, we generate midpoints using a generative A\* path planning method, following Wang et al. [55].

The *dynamic* setting involves five pairs of object and human starting locations, accompanied by a trajectory leading towards the object. Each method is tasked with creating five unique motion variations that both approach and interact with the object, conforming to a designated action type.

### A.2. Implementation Details

Our motion generation model utilizes a DDPM architecture with a linear variance scheduler, conditioned on scene and action embeddings. Following Guo et al. [13], we implement a Transformer encoder as the UNet structure with an embedding dimensionality of 512 for projecting body joint locations into a high-dimensional space. The Transformer consists of 6 layers, 16 heads in the multi-head attention, and uses a 0.1 dropout rate. The intermediate feedforward network’s dimensionality is set to 1024. Both scene and action encoders are Transformer-based with 6 layers and 16 heads, producing 512-dimensional embeddings. These embeddings are added to the first token of the motion generation transformer, which contains timestep information. The Huber loss is used to gauge the accuracy of predicted noise.

For converting joint locations to SMPL-X parameterized meshes, a pre-trained 4-layer MLP predicts coarse SMPL-X parameters, with a 6D representation for rotation [66]. The MLP inputs three consecutive frames and outputs parameters for the middle frame. Edge cases use duplicated middle frames for input. An optimization process refines body poses using gradient descent on the L2 error between joint locations from the model and predicted SMPL-X parameters, ensuring accurate body pose representation.

Training with the Adam optimizer (learning rate  $1e-4$ ,

batch size 512) on four NVIDIA A800 GPUs, our method takes 48 hours to train for 500k steps on TRUMANS.

### A.3. Adaption of Baselines

We adapted baseline methods for a fair comparison in our autoregressive long-term motion generation framework. For Wang et al. [53], two milestone poses are set as the transition and subgoal at the start and end of each episode, with in-between motions generated using their original cVAE model. Methods like Huang et al. [21] and Karunratanakul et al. [23], not initially designed for long-term synthesis, were modified to incorporate  $M_{goal}$  and  $M_{trans}$  in the sampling stage, maintaining their original sampling strategies. In *dynamic* experiments involving dynamic objects, we adapted models like GOAL [48] to encompass both reaching and grasping actions, while preserving their original training pipeline.

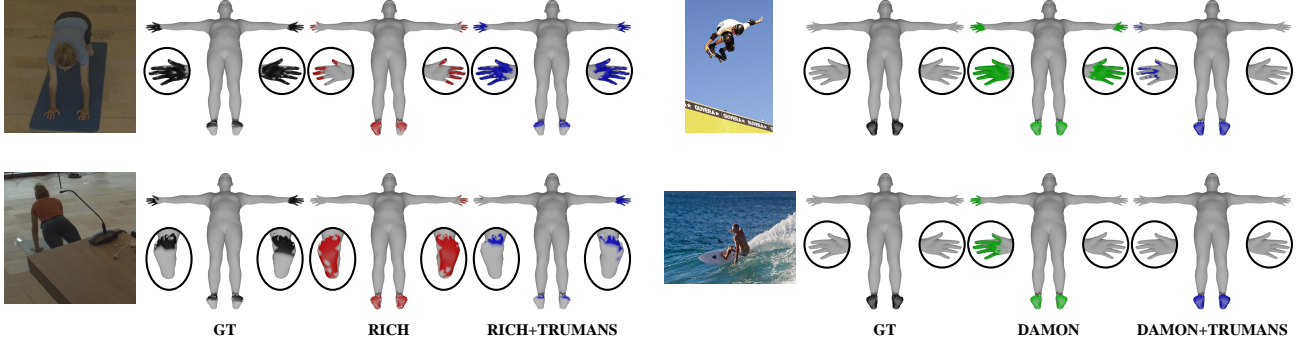
### A.4. Human Study

We conducted human studies with 31 participants (17 males and 14 females) in a controlled environment, ensuring no communication among them. For each baseline and ablated version of our method, we generated 4 human motions aligned with the respective test settings (dynamic or non-dynamic). These motions were applied to the SMPL-X body mesh and rendered in the same scene from TRUMANS-test for horizontal comparison. Alongside these synthesized motions, one MoCap motion from TRUMANS-test in the same scene was also included, with careful rendering to minimize visual obstructions.

To qualitatively compare our method with SAMP [17], DIMOS [65], LAMA [25], and Wang et al. [55], we replicated their demonstrations using our model trained on TRUMANS. This involved setting similar subgoals to duplicate the trajectories. Participants were shown side-by-side comparisons of motions synthesized by our method and the baseline methods’ demos and then asked to choose the more natural-looking output. The frequency of our method being preferred is reflected in the Success Rate of Discrimination (SucRateDis) reported in Tab. A1.

Table A1. **Human study results of comparisons between our method with recent work.** The Success Rate of Discrimination (SucRateDis), indicating the frequency at which our method is selected as the superior one, is reported.

Method	Success Rate of Discrimination (%)
SAMP [17]	100
Wang et al. [55]	100
LAMA [25]	80.6
DIMOS [65]	64.5



(a) BSTRO [20] trained on RICH [20] combined with TRUMANS. (b) DECO [50] trained on DAMON [50] combined with TRUMANS.

Figure A1. Additional qualitative results of 3D contact estimation.

## A.5. Image-based Tasks

This section details additional qualitative and quantitative results for image-based tasks using the rendered images and annotations from TRUMANS.

**3D Human Mesh Estimation** To assess the impact of integrating TRUMANS into training, we use two additional methods, I2L [34] and SGRE [52], on the 3DPW test set. These methods are trained either solely on 3DPW or on a combination of 3DPW and TRUMANS at varying ratios. As indicated in Tabs. A2 and A3, incorporating our synthetic data with the real training set markedly enhances performance.

Table A2. Performance of I2L [34] in 3D human mesh estimation trained on 3DPW [51] combined with TRUMANS in different ratios.

Training Data	MPVE↓	MPJPE↓	PA-MPJPE↓
3DPW [51]	186.9	160.4	90.2
3DPW+T (2:1)	133.2	116.5	69.1
3DPW+T (1:1)	<b>126.1</b>	<b>110.2</b>	<b>66.2</b>

Table A3. Performance of SGRE [52] in 3D human mesh estimation trained on 3DPW [51] combined with TRUMANS in different ratios.

Training Data	MPVE↓	MPJPE↓	PA-MPJPE↓
3DPW [51]	257.0	223.0	110.6
3DPW+T (2:1)	240.6	207.2	113.5
3DPW+T (1:1)	<b>138.0</b>	<b>117.5</b>	<b>80.3</b>

**3D Contact Estimation** Qualitative results of baseline methods trained with and without TRUMANS are presented in Fig. A1. These results demonstrate that the incorporation of TRUMANS in training enhances the precision of contact prediction.

## B. Additional Details of TRUMANS

### B.1. Additional Details of Dataset Comparison

In our dataset comparison presented in Tab. 1, we have categorized similar objects into common types for a more equitable comparison. For SAMP [17], their seven reported objects, including various chairs and a table, are grouped into three types: “sofa,” “chair,” and “table.” BEHAVE [2], with a list of 20 distinct items, is classified into 14 object types, consolidating similar items like chairs and tables. Similarly, iReplica [15]’s report of 9 objects is condensed into 5 classes.

Additionally, for iReplica, we have combined data from their two datasets, EgoHOI and H-contact, for simplicity. EgoHOI contributes 0.25 hours of HSI data with ego-view and multiview RGBD videos, while H-contact adds 0.5 hours of HSI data featuring per-frame hand-object contact.

### B.2. Dataset Splits

TRUMANS is divided into training, validation, and test sets, with scenes 1 to 70 for training, 71 to 80 for validation, and 91 to 100 for testing. This distribution creates a split ratio of approximately 7:1:2 across all data frames for the respective sets.

### B.3. Object Types

TRUMANS includes 20 types of objects commonly found in indoor scenes, categorized as either [a] (articulated) or [r] (rigid). The list with descriptions is as follows:

- Articulated chair: [a], including gaming and office chairs.
- Rigid chair: [r], encompasses chairs with/without armrests and stools.
- Table: [r], available in round and square shapes.
- Sofa: [r], varieties like single-seaters and couches.
- Bed: [r].
- Book: [a].
- Pen: [r].

- Phone: [r].
- Mouse: [r].
- Keyboard: [r].
- Handbag: [r].
- Vase: [r].
- Cup: [r].
- Bottle: [r].
- Laptop: [a].
- Oven: [a].
- Drawer: [a].
- Cabinet: [a].
- Microwave: [a].
- Door: [a].

## B.4. Capture Pipeline

**Aligning virtual and real environments** The alignment between virtual and real environments is crucial to ensure the plausibility of actions projected into the virtual world. Our process starts with manually selecting scenes and objects from the 3D-FRONT [10] dataset and BlenderKit [6], prioritizing scenes with numerous interactable objects that fit within our motion capture area. Manual alignment is performed to match these virtual scenes with real-world counterparts. For example, a digital sofa may be replicated with a real sofa or chairs arranged to mimic its shape. When digital and physical shapes do not match perfectly, we modify the digital asset, such as scaling objects or editing meshes, or adjust our real-world setups, like placing a mat on a chair to simulate a higher seat.

To align digital characters with our real actors, we start by exporting a human armature matching the actor’s bone lengths using the VICON Shogun system [41]. Then, in Character Creator 4 [38], we adjust sliders to create digital humans mirroring the actors’ bone lengths. These digital characters are re-imported into Shogun for real-time IK re-targeting, ensuring accurate character poses for our digital humans.

**Object placeholders** Our dataset addresses the limitations of previous datasets related to object visibility and complex HSIs in clustered scenes. We use an optical MoCap system with object placeholders designed for light transmission, enabling accurate capture even in complex HSIs. For example, an actor seated behind a transparent acrylic table allows for precise leg tracking.

**Real-time data quality inspection** Data recording is monitored in real-time to ensure quality. Inspectors watch the digital world and human avatars on screens during capture, filtering out obvious errors like untracked markers or jittering in IK solving. This real-time inspection aids in maintaining high data quality.

## B.5. Motion Augmentation Implementation Details

This section delves into the specifics of our motion augmentation pipeline, illustrating the need for our refinement processes with examples and establishing theoretical bounds for the smoothness of IK target trajectories.

In our augmentation process, we first identify contact points between human joints and object meshes in a given motion sequence. For example, if joint  $J_1$  is in contact with an object mesh at point  $\mathbf{v}_m$  at time  $T_1$ , and we subsequently alter the object’s shape or replace it, the new corresponding point becomes  $\mathbf{v}'_m$ . To preserve the interaction, the joint’s target location  $\mathbf{l}'$  must be adjusted accordingly to maintain the contact point:

$$\mathbf{l}'_{T_1} - \mathbf{v}'_m = \mathbf{l}_{T_1} - \mathbf{v}_m, \quad (\text{A1})$$

such that

$$\mathbf{l}'_{T_1} = \mathbf{l}_{T_1} + \mathbf{v}'_m - \mathbf{v}_m. \quad (\text{A2})$$

The offset  $\mathbf{v}_1 = \mathbf{v}'_m - \mathbf{v}_m$  represents the change in joint position resulting from the object’s shape variance. To address the abrupt trajectory change, we implement a smoothing process for the pose trajectory. Within a defined proximity window  $W$ , we apply this offset with a linearly decreasing norm to ensure a smoother transition in the joint’s trajectory:

$$\mathbf{v}_{1t} = \left(1 - \frac{|t - T_1|}{W}\right) \mathbf{v}_1. \quad (\text{A3})$$

This offset application ensures a seamless blend from the original trajectory to the augmented one. To validate the smoothness of this trajectory, we establish a bound on the difference in bone target positions between consecutive frames. The notation  $|\cdot|$  is used to denote the absolute value for scalars and the 2-norm for vectors, aiding in quantifying these trajectory modifications.

$$\begin{aligned} |\mathbf{l}'_{t+1} - \mathbf{l}'_t| &= |\mathbf{l}_{t+1} + \mathbf{v}_{t+1} - \mathbf{l}_t - \mathbf{v}_t| \\ &\leq |\mathbf{l}_{t+1} - \mathbf{l}_t| + |\mathbf{v}_{t+1} - \mathbf{v}_t| \\ &\leq |\mathbf{l}_{t+1} - \mathbf{l}_t| + \frac{|\mathbf{v}|}{W}. \end{aligned} \quad (\text{A4})$$

When the smoothing window length  $W$  is set appropriately, such as 30 in our implementation, and the norms of the offset vectors are within practical limits (dictated by reasonable object variation scales and stable CCD IK results), the updated IK target trajectories exhibit sufficient smoothness.

However, this IK solving approach is limited to simpler scenarios involving interaction with a single object. In more complex situations, typical of our dataset, we encounter additional challenges. For instance, consider an offset  $\mathbf{v}_1$  introduced at time  $T_1$  for bone  $J_1$  due to object augmentation, with this contact ending at  $T_2$ . At  $T_2$ , another bone  $J_2$  enters into contact, necessitating a new IK target. Post-IK process, bones without specific targets, including  $J_1$ , may shift from

their original positions. We denote the deviation for bone  $J_1$  as  $\mathbf{v}_2$ . To mitigate this deviation, we employ a similar offsetting technique by blending  $\mathbf{v}_1$  and  $\mathbf{v}_2$ . For each time  $t$  within this window of length  $W$ , the bone is assigned an offset vector that is a weighted mean of these two offsets, calculated using vector norms as weights:

$$\text{offset} = \frac{(|\mathbf{v}_{1t}|\mathbf{v}_{1t} + |\mathbf{v}_{2t}|\mathbf{v}_{2t})}{|\mathbf{v}_{1t}| + |\mathbf{v}_{2t}|}, \quad (\text{A5})$$

where

$$\begin{aligned} \mathbf{v}_{1t} &= \left(1 - \frac{|t - T_1|}{W}\right)\mathbf{v}_1, \\ \mathbf{v}_{2t} &= \left(1 - \frac{|t - T_2|}{W}\right)\mathbf{v}_2. \end{aligned} \quad (\text{A6})$$

By integrating these two stages of linear blending of offset vectors, we achieve smooth trajectories for IK targets. As outlined earlier, akin to the approach in Eq. (A4), we analyze the joint target differences between consecutive frames to further substantiate the smoothness of our approach:

$$\begin{aligned} & |I'_{t+1} - I'_t| \\ = & \left| \mathbf{l}_{t+1} - \mathbf{l}_t \right. \\ & + \frac{|\mathbf{v}_{1,t+1}|\mathbf{v}_{1,t+1} + |\mathbf{v}_{2,t+1}|\mathbf{v}_{2,t+1}}{|\mathbf{v}_{1,t+1}| + |\mathbf{v}_{2,t+1}|} \\ & \left. - \frac{|\mathbf{v}_{1,t}|\mathbf{v}_{1,t} + |\mathbf{v}_{2,t}|\mathbf{v}_{2,t}}{|\mathbf{v}_{1,t}| + |\mathbf{v}_{2,t}|} \right| \end{aligned} \quad (\text{A7})$$

$$\leq \left| \mathbf{l}_{t+1} - \mathbf{l}_t \right| \quad (\text{A8})$$

$$+ \left| \frac{|\mathbf{v}_{1,t+1}|\mathbf{v}_{1,t+1}}{|\mathbf{v}_{1,t+1}| + |\mathbf{v}_{2,t+1}|} - \frac{|\mathbf{v}_{1,t}|\mathbf{v}_{1,t}}{|\mathbf{v}_{1,t}| + |\mathbf{v}_{2,t}|} \right| \quad (\text{A9})$$

$$+ \left| \frac{|\mathbf{v}_{2,t+1}|\mathbf{v}_{2,t+1}}{|\mathbf{v}_{1,t+1}| + |\mathbf{v}_{2,t+1}|} - \frac{|\mathbf{v}_{2,t}|\mathbf{v}_{2,t}}{|\mathbf{v}_{1,t}| + |\mathbf{v}_{2,t}|} \right|, \quad (\text{A10})$$

where

$$|\mathbf{v}_{i,t}| = \left(1 - \frac{|t - T_i|}{W}\right)|\mathbf{v}_i|. \quad (\text{A11})$$

Thus we have

$$\begin{aligned} & \left| \frac{|\mathbf{v}_{1,t+1}|\mathbf{v}_{1,t+1}}{|\mathbf{v}_{1,t+1}| + |\mathbf{v}_{2,t+1}|} - \frac{|\mathbf{v}_{1,t}|\mathbf{v}_{1,t}}{|\mathbf{v}_{1,t}| + |\mathbf{v}_{2,t}|} \right| \\ = & \frac{|\mathbf{v}_{1,t+1}|\mathbf{v}_{1,t}(|\mathbf{v}_{1,t+1} - \mathbf{v}_{1,t}|)}{(|\mathbf{v}_{1,t+1}| + |\mathbf{v}_{2,t+1}|)(|\mathbf{v}_{1,t}| + |\mathbf{v}_{2,t}|)} \\ & + \frac{(|\mathbf{v}_{1,t+1}|\mathbf{v}_{2,t}|\mathbf{v}_{1,t+1} - |\mathbf{v}_{2,t+1}|\mathbf{v}_{1,t}|\mathbf{v}_{1,t})}{(|\mathbf{v}_{1,t+1}| + |\mathbf{v}_{2,t+1}|)(|\mathbf{v}_{1,t}| + |\mathbf{v}_{2,t}|)} \\ \leq & \frac{|\mathbf{v}_{1,t+1}|\mathbf{v}_{1,t}|\mathbf{v}_{1,t+1} - \mathbf{v}_{1,t}|}{(|\mathbf{v}_{1,t+1}| + |\mathbf{v}_{2,t+1}|)(|\mathbf{v}_{1,t}| + |\mathbf{v}_{2,t}|)} \\ & + \frac{|\mathbf{v}_{1,t+1}|\mathbf{v}_{2,t}|\mathbf{v}_{1,t+1} - |\mathbf{v}_{2,t+1}|\mathbf{v}_{1,t}|\mathbf{v}_{1,t}|}{(|\mathbf{v}_{1,t+1}| + |\mathbf{v}_{2,t+1}|)(|\mathbf{v}_{1,t}| + |\mathbf{v}_{2,t}|)} \\ & + \frac{|\mathbf{v}_{1,t+1}|\mathbf{v}_{2,t}|\mathbf{v}_{1,t} - |\mathbf{v}_{1,t}|\mathbf{v}_{2,t}|\mathbf{v}_{1,t}|}{(|\mathbf{v}_{1,t+1}| + |\mathbf{v}_{2,t+1}|)(|\mathbf{v}_{1,t}| + |\mathbf{v}_{2,t}|)} \\ & + \frac{|\mathbf{v}_{1,t}|\mathbf{v}_{2,t}|\mathbf{v}_{1,t} - |\mathbf{v}_{1,t}|\mathbf{v}_{2,t+1}|\mathbf{v}_{1,t}|}{(|\mathbf{v}_{1,t+1}| + |\mathbf{v}_{2,t+1}|)(|\mathbf{v}_{1,t}| + |\mathbf{v}_{2,t}|)} \\ < & \frac{|\mathbf{v}_{1,t+1}|\mathbf{v}_{1,t}|\mathbf{v}_{1,t+1} - \mathbf{v}_{1,t}|}{|\mathbf{v}_{1,t+1}|\mathbf{v}_{1,t}|} \\ & + \frac{|\mathbf{v}_{1,t+1}|\mathbf{v}_{2,t}|\mathbf{v}_{1,t+1} - \mathbf{v}_{1,t}|}{|\mathbf{v}_{1,t+1}|\mathbf{v}_{2,t}|} \\ & + \frac{|\mathbf{v}_{1,t}|\mathbf{v}_{2,t}(|\mathbf{v}_{1,t+1}| - |\mathbf{v}_{1,t}|)}{|\mathbf{v}_{1,t}|\mathbf{v}_{2,t}|} \\ & + \frac{|\mathbf{v}_{1,t}|\mathbf{v}_{1,t}(|\mathbf{v}_{2,t+1}| - |\mathbf{v}_{2,t}|)}{|\mathbf{v}_{2,t+1}|\mathbf{v}_{2,t}|} \\ = & \frac{3|\mathbf{v}_1| + |\mathbf{v}_2|}{W}. \end{aligned} \quad (\text{A12})$$

With the scaling of Eqs. (A9) and (A10) now aligned, we substitute these elements to establish a theoretical bound:

$$\begin{aligned} & |I'_{t+1} - I'_t| \\ < & |\mathbf{l}_{t+1} - \mathbf{l}_t| + \frac{3|\mathbf{v}_1| + |\mathbf{v}_2|}{W} + \frac{|\mathbf{v}_1| + 3|\mathbf{v}_2|}{W} \\ = & |\mathbf{l}_{t+1} - \mathbf{l}_t| + \frac{4}{W}(|\mathbf{v}_1| + |\mathbf{v}_2|), \end{aligned} \quad (\text{A13})$$

which ensures that our target trajectories are smooth.

## B.6. Annotations

**Human Motion** The motion data captured from the VICON system, initially in the form of pose sequences of a custom armature, is converted into the SMPL-X format [36] using a vertex-to-vertex optimization method. This ensures accurate and smooth SMPL-X representations; please refer to Fig. A2 for examples. The conversion process involves the following steps:

1. Vertices on the SMPL-X mesh are manually selected and paired with the closest vertices on our custom mesh.
2. A loss function, defined as the Mean Squared Error (MSE) between paired vertex locations, is minimized using the Adam optimizer to refine SMPL-X parameters until convergence.

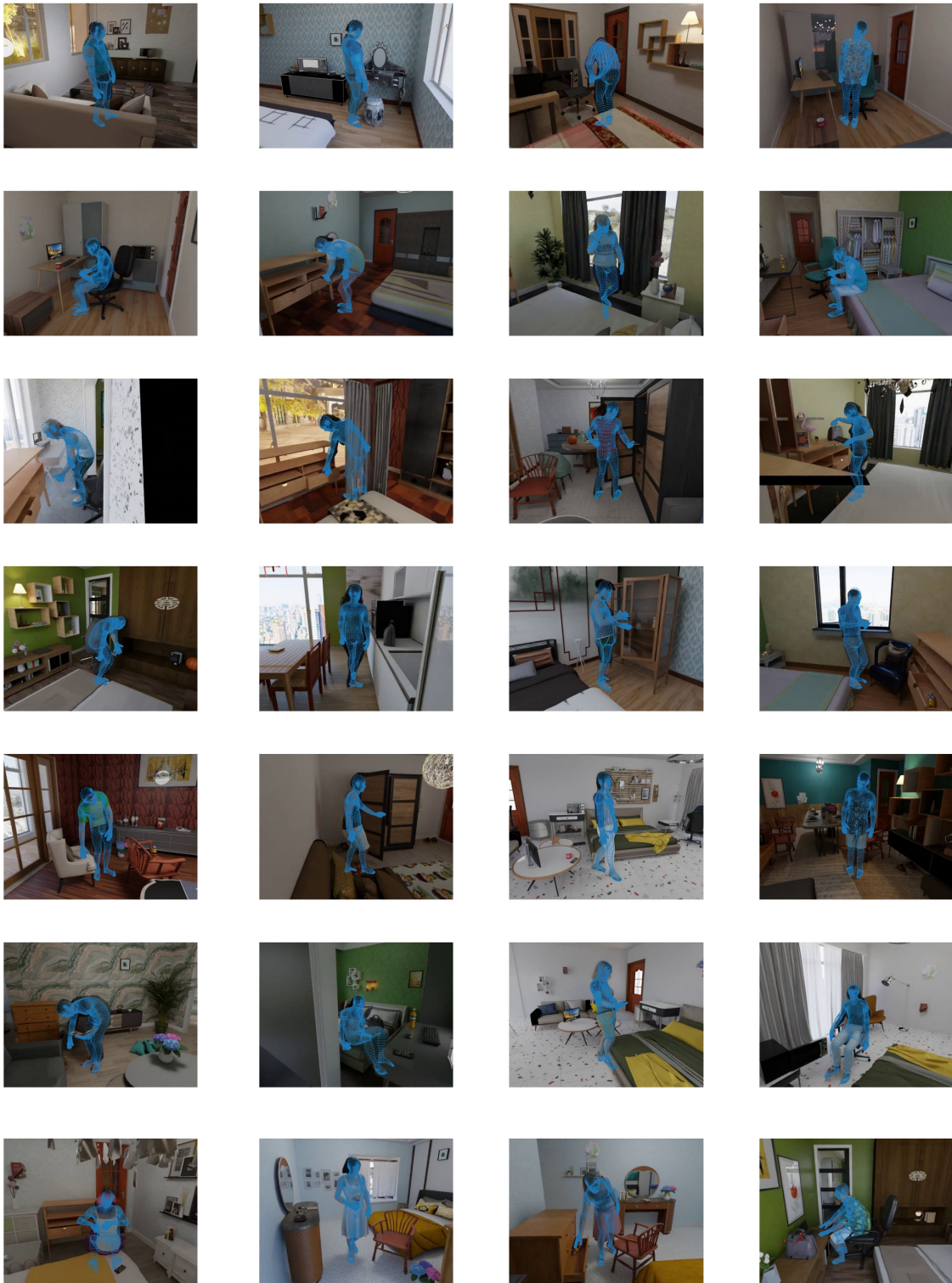


Figure A2. Examples of SMPL-X annotations in TRUMANS.



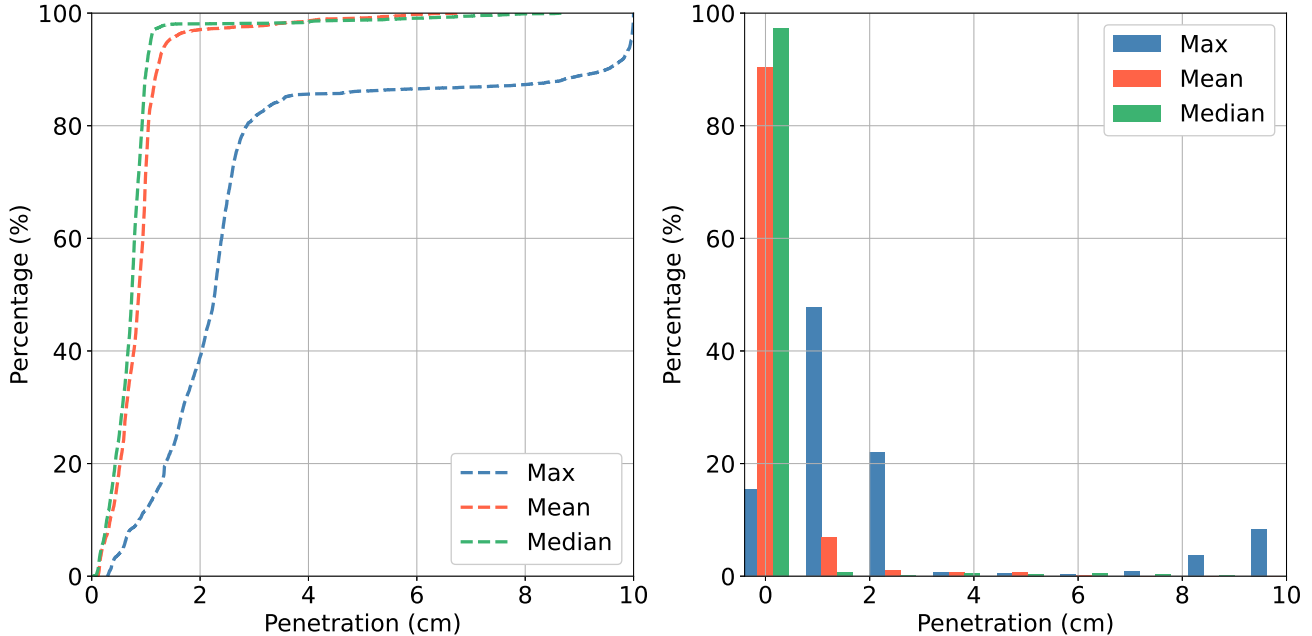


Figure A3. **Penetration statistics in TRUMANS.** This analysis covers maximum, mean, and median penetration distances per vertex per frame. The left graph displays the proportion of frames with penetration below various thresholds (X-axis), and the right bar plot categorizes frames by specific penetration distance ranges (X-axis). Notably, in more than 95% of frames, both the mean and median penetration distances stay below 2cm.

3. Inaccuracies in the SMPL-X mesh are manually corrected by adjusting bone rotations.

After aligning the SMPL-X mesh with our custom mesh, we record the mapping for future use.

In the second phase, we enhance the custom motion data by adding interpolated frames from the T-pose to the start pose. This ensures a smooth transition for each bone in the sequence.

Finally, we optimize the SMPL-X parameters frame by frame, starting from the pose established in the first phase. We first refine the body shape parameters and then adjust the pose parameters, including global translation. The optimization of each frame starts from the pose of the previous frame and continues until convergence. This method relies on minimal mesh changes between frames, supported by our high-quality motion data. A typical MSE value at convergence ranges between  $5e-5$  and  $1e-4$ , indicating an average point distance of less than 1cm.

**Contact** Following the fitting of SMPL-X meshes, we compute per-vertex contact annotations. The contact for each human mesh vertex is determined based on its proximity and orientation relative to scene or object meshes. A vertex is deemed in contact if it fulfills either of the following conditions: (i) it resides inside an object’s mesh, or (ii) while outside an object’s mesh, it is within a specified threshold distance and the angle between its normal and the vector pointing towards the object is under 60 degrees. The latter

criterion is particularly vital for accurate contact annotation, as it prevents misidentification in scenarios like a hand holding a bottle. Penetration statistics, as detailed in Fig. A3, reveal that in over 95% of the frames, both the mean and median penetration distances remain below 2cm. For examples of contact annotation, please refer to Fig. A4.

**Objects** In TRUMANS, we include the watertight mesh for all objects and their 6D poses for each frame. For articulated objects, part-level annotations are provided along with the URDF (Unified Robot Description Format) files to represent their kinematic structure.

**Actions** For every sequence within TRUMANS, multi-hot action labels are assigned on a frame-by-frame basis. This approach allows for the representation of multiple concurrent actions in a single frame, with each action identified by a distinct label.

## B.7. Video Rendering

To enhance video diversity and accurately capture HSI details, we developed an adaptive camera tracking algorithm that maintains footage smoothness and consistency. The camera, set at a constant height of 1.4 meters, moves within a 2-meter radius around the human, horizontally oriented towards the body. The camera’s pose is exclusively determined by its rotation around the z-axis in the human coordinate system.

We set keyframes at intervals of 30 frames. For each keyframe, 20 camera proposals adhering to the constraints are pre-defined and evenly distributed around the ring. To identify active hand interactions, we calculate the minimum distance between hand joints and dynamic objects. If this distance exceeds 20 centimeters, we default to tracking the right hand. For the identified interacting hand, rays emitted from each camera proposal towards the hand joints help measure visibility. A joint is considered visible to a camera if the intersection point with the scene's mesh can be projected onto the camera's imaging plane, and the distance to the joint is less than 10 centimeters. The number of visible joints determines whether a camera effectively captures the interaction. The visibility threshold for different keyframes is dynamically adjusted to ensure at least one camera captures the interaction, except when all joints are invisible.

After assessing the interaction capture capability of the 20 cameras at each keyframe, dynamic programming is used to select the optimal keyframe camera pose sequence that minimizes total rotation and maximizes interaction coverage. "Camera pose" here specifically refers to rotation about the z-axis. Camera poses for frames between keyframes are interpolated using cubic spline interpolation of the rotation angles.



Figure A4. Examples of contact annotations in TRUMANS.