

EmbodMocap: In-the-Wild 4D Human-Scene Reconstruction for Embodied Agents

Wenjia Wang¹ Liang Pan¹ Huaijin Pi¹ Yuke Lou¹
 Xuqian Ren³ Yifan Wu¹ Zhouyingcheng Liao¹ Lei Yang² Taku Komura¹

¹The University of Hong Kong ²The Chinese University of Hong Kong ³Tampere University

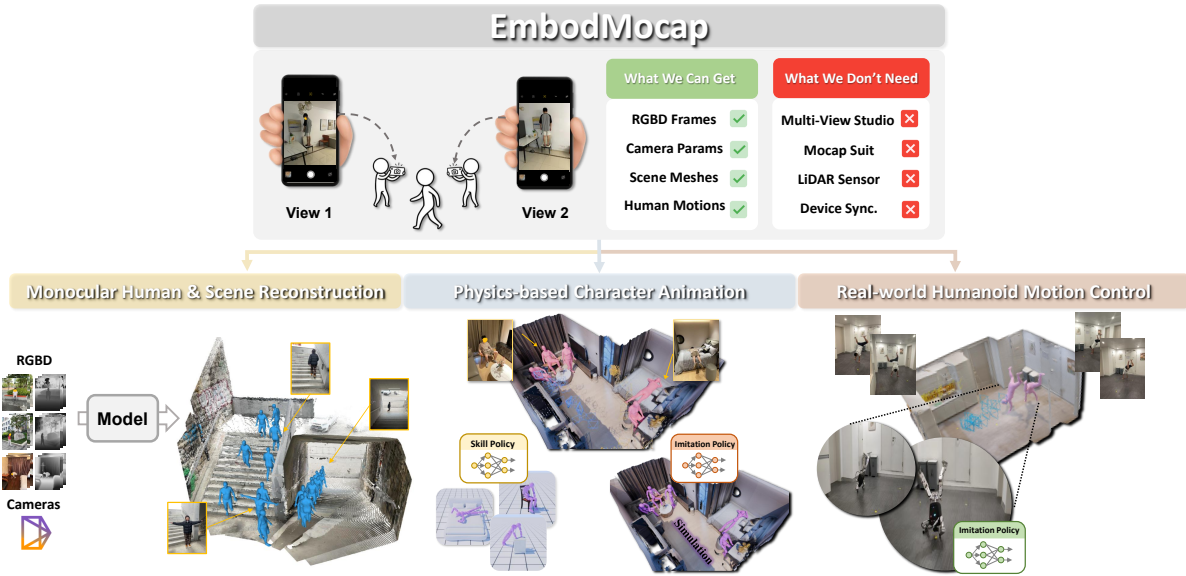


Figure 1. Introducing **EmbodMocap**, a portable and low-cost system for simultaneous 4D human and scene reconstruction, deployable anywhere using two moving iPhones. The dataset captured by EmbodMocap benefits three crucial embodied AI tasks: monocular human & scene reconstruction, physics-based character animation, and real-world humanoid motion control.

Abstract

Human behaviors in the real world naturally encode rich, long-term contextual information that can be leveraged to train embodied agents for perception, understanding, and acting. However, existing capture systems typically rely on costly studio setups and wearable devices, limiting the large-scale collection of scene-conditioned human motion data in the wild. To address this, we propose *EmbodMocap*, a portable and affordable data collection pipeline using two moving iPhones. Our key idea is to jointly calibrate dual RGB-D sequences to reconstruct both humans and scenes within a unified metric world coordinate frame. The proposed method allows metric-scale and scene-consistent capture in everyday environments without static cameras or markers, bridging human motion and scene geometry seamlessly. Based on the collected data, we empower three embodied AI tasks: monocular human-scene-

reconstruction, where we fine-tune on feedforward models that output metric-scale, world-space aligned humans and scenes; physics-based character animation, where we prove our data could be used to scale human-object interaction skills and scene-aware motion tracking; and robot motion control, where we train a humanoid robot via sim-to-real RL to replicate human motions depicted in videos. Experimental results validate the effectiveness of our pipeline and its contributions towards advancing embodied AI research.

1. Introduction

Embodied Artificial Intelligence (Embodied AI) aims to build agents that can perceive, understand, and act within real-world environments. Progress in this field relies on datasets that capture both human motion and the surrounding 3D scene, enabling physically grounded perception and action learning. Such scene-aware data allows modeling

of realistic human–scene interactions, simulation of lifelike behaviors, and training of humanoids to operate seamlessly in complex environments. They serve as a foundation for advancing embodied reasoning and control across robotics, virtual reality, and computer vision.

However, collecting high-quality human–scene data remains difficult. Precise 3D motion and scene geometry cannot be automatically obtained from internet videos due to occlusions and depth ambiguity. Existing capture systems that provide high-quality human–scene data typically rely on multi-view camera rigs [11, 74], wearable motion suits [22, 35], or LiDAR scanners [6, 19], which are costly, complex, and limited to controlled studio environments. These constraints hinder scalable and scene-aware data acquisition, limiting the ability of embodied AI models to learn from natural human behavior in diverse indoor and outdoor environments.

In this paper, we propose EmbodMocap, an efficient and affordable framework for capturing metrically accurate 4D human and scene using only two iPhones. Our key idea is to jointly calibrate and optimize dual RGB-D inputs to reconstruct both humans and scenes within a unified world coordinate frame. Specifically, we first reconstruct the static scene from a single RGB-D sequence to define the world scale, then capture synchronized dual-view RGB-D videos of human motion, and finally perform geometric alignment and motion optimization to recover world-anchored human poses. In contrast to existing systems that rely on multi-camera rigs or wearable sensors, our approach achieves high-quality, scene-consistent reconstruction using only moving consumer devices. This design enables scalable, in-the-wild data collection that preserves precise human motion and authentic scene context, supporting realistic human–scene interaction modeling for embodied AI research.

Based on the data collected with EmbodMocap, we demonstrate the reliability and versatility of our capture pipeline through three representative applications. The first application verifies geometric consistency, where we fine-tune reconstruction models to jointly recover humans and scenes in world coordinates. The second validates physical realism, showing that the captured motions enable scalable training of physics-based character skills and scene-aware motion tracking. The third demonstrates embodied transferability, where our data support humanoid robot training through a sim-to-real motion tracking framework [26, 44]. These results highlight that EmbodMocap enables scalable and physically grounded data acquisition for embodied AI.

In summary, our contributions can be summarized as follows:

- We introduce EmbodMocap, a portable and affordable data collection pipeline that produces high-quality multi-modal data for embodied AI applications.

- We validate our capture pipeline’s effectiveness across three key embodied AI tasks: monocular human-scene reconstruction, physics-based character animation, and real-world humanoid motion control.
- We provide a scalable and accessible solution that lowers the barrier for embodied AI research, opening new possibilities for real-world applications and further advancements in the field. All the codes and datasets will be open-sourced.

2. Related Work

Datasets for 4D Human & Scene Capture. Early motion datasets, such as AMASS [10, 36], focus on pure human motion, unifying multiple motion capture sources into a large-scale repository. While invaluable for studying human motion, these datasets lack the 3D scene context essential for understanding human–scene interactions. Recent 4D datasets, like PROX [11], RICH [19], and Ego-Body [74], combine scanned 3D scenes with motion capture using multi-view camera systems, while EMDb [22] and SPLOPER4D [6], employ IMUs or electromagnetic sensors for motion recording in large-scale environments. Nymeria [35] extends this further with Project Aria glasses and optical marker-based systems for wide-area motion capture. However, these approaches face notable limitations: marker-based and multi-camera systems are expensive and restricted to small studio environments, while IMU and EM-based methods, though more flexible, require extensive manual alignment and post-processing to synchronize motion with 3D scenes. And the wearable devices will influence the human appearance in RGB images. In contrast, our approach uses minimal equipment, operates in diverse environments without static camera setups, and avoids wearable devices, preserving the naturalness of RGB images for authentic human–scene interaction capture. Table 1 compares these datasets.

Monocular Human & Scene Reconstruction. Early works [4, 8, 21, 24, 42] on RGB-based human mesh recovery focus on reconstructing 3D pose and shape but often ignore scene context [60] or camera information [25, 63], leading to inconsistencies under camera motion. Recent methods address this by combining motion cues [73], SLAM or visual odometry [55, 65, 72], and human motion priors [54, 73] to recover global trajectories in world coordinates.

Emerging models move toward jointly reconstructing humans and 3D scenes with spatial intelligence models [61, 62]. For example, HSFm [38] combines Dust3R [62] with multi-view correspondence to jointly recover human meshes, scene point clouds, and camera parameters from multi-cameras. HAMSt3R [49] integrates DensePose [9] and multi-view scene reconstruction in one model, with an optimization to get human poses, while JOSH [29]

Table 1. Comparison of 4D Human & Scene datasets based on different features.

Datasets	Publication	Device					Outcome		
		Mocap Suit	Scanner	Static Cam.	Dyna. Cam.	Total Cost(\$)	Mesh	Dyna.Anno.	Outdoor
PROX [11]	ICCV2019	-	Structure Sensor	Kinetic-One	-	2K	✓	✗	✗
RICH [19]	CVPR 2022	-	Leica RTC360	6-8×Cameras	1×Camera	20K+	✓	✓	✓
EgoBody [74]	ECCV2022	-	1×iPhone	5×Azure Kinect	Hololens2	9K	✓	✓	✗
SLOPER4D [6]	CVPR2023	Noitom PN+NUC11	Ouster-os1 LiDAR	-	DJI-Action2+TLS	20K	✓	✓	✓
EMDB [22]	ICCV 2023	EM Sensors	-	-	1×iPhone	15K	✗	✓	✓
Nymeria [35]	ECCV2024	2×XSens+Aria Wistband	-	-	2×Project Aria	60K+	✗	✓	✓
EmbodMocap	-	-	1×iPhone	-	2×iPhone	1K	✓	✓	✓

uses MAST3R-SLAM [39] and joint optimization to achieve globally consistent 4D human-scene reconstructions. This trend emphasizes the simultaneous prediction of human motion and scene geometry, which further requires multi-model data pairs with high-quality annotations. In our paper, we propose a monocular human & scene reconstruction pipeline combined with 2 feedforward models, and fine-tuned it on our proposed dataset to prove the efficiency of our paired data.

Training Humanoid from Video Data. Recent advances in physics-based animation and reinforcement learning enable humanoid agents to perform realistic and physically consistent motions using control policies learned from marker-based motion capture data. These methods have shown strong realism in tasks like motion tracking [32, 44], locomotion [33, 45, 46], and human-scene interaction [41, 64], and have been extended to real-world applications in motion tracking [15, 17, 20], locomotion [16], and scene interaction [3, 14]. However, marker-based methods require dedicated studios, expensive hardware, and extensive manual effort, making them costly and hard to scale. Adapting captured motions to new scenes or robot morphologies also demands complex retargeting and re-simulation. To address this, recent works like VideoMimic [2], ASAP [17], and HDMI [67] train humanoid control directly from in-the-wild video data. By using monocular motion capture methods such as TRAM [65] and GVHMR [54], they estimate human motion from videos and retarget it to virtual humanoids for training in physical simulators. This video-driven paradigm leverages diverse real-world data but struggles with capturing complex skills or scene geometries due to occlusion and depth ambiguities. In this paper, we propose a method for high-precision human motion and scene reconstruction that overcomes these limitations.

3. Proposed Capture System

We aim to capture metrically accurate human motion and scene geometry using only two iPhones. As shown in Fig. 2, our capture process consists of four sequential stages that progressively reconstruct and align the scene, cameras, and human motion within a unified world coordinate frame. We first reconstruct a metrically accurate static scene and establish the world reference using a single iPhone RGB-

D sequence (Sec. 3.1). Then, we use two synchronized iPhones to record dual-view RGB-D videos of human motion and extract per-frame camera poses and human priors with off-the-shelf perception models (Sec. 3.2). Next, we align the dual-view camera trajectories to the reconstructed scene through a combination of COLMAP registration and multi-view geometric optimization (Sec. 3.3). Finally, we refine the SMPL parameters by triangulating dual-view 2D keypoints into 3D space and optimizing human poses and translations in the world coordinate system (Sec. 3.4).

3.1. Stage I: Scene Reconstruction

In this stage, we aim to reconstruct a metrically accurate, Z-up scene mesh that serves as the reference world coordinate system. We first use a single iPhone to capture an RGB-D video of the scene, along with synchronized IMU data. The recorded data are processed by the SpectacularAI SDK (SAI) [1], which automatically selects keyframes according to the accumulated camera translation and estimates corresponding camera parameters ($K_s, R_{s,n}, T_{s,n}$) in Z-up world coordinates with metric scale. These trajectories establish a consistent world frame for all subsequent stages. Based on the recovered poses, we refine the iPhone LiDAR depth maps using PromptDA [27], unproject them into 3D space, and integrate the point clouds through TSDF fusion [5] to obtain a dense and metrically accurate global mesh \mathcal{M}_g . Note that the depth maps are truncated based on a threshold determined by the effective range of the iPhone’s depth sensor. Specifically, we use a threshold of 3.5m for indoor scenes and 5m for outdoor scenes. We further apply lightweight post-processing such as outlier removal and small-component filtering to clean the mesh. Finally, we extract SIFT features from the same SAI keyframes and run COLMAP [51] with fixed camera parameters to build a sparse structure database. This database preserves the metric scale and serves as a reference for registering dual-view sequences in later stages.

3.2. Stage II: Sequence Processing

After reconstructing the static scene in Stage I, we proceed to capture and process dual-view human motion sequences within the same environment. In this stage, we use two iPhones to record synchronized RGB-D videos of a per-

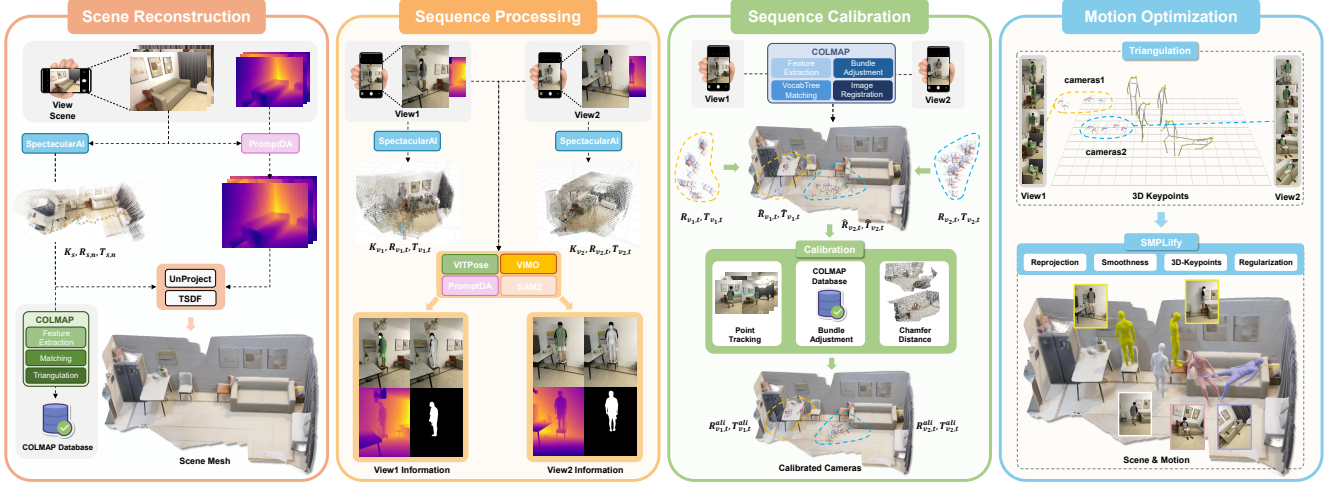


Figure 2. EmbodMocap: We propose an affordable dataset capture and processing system. From left to right, the four stages (Stage-I to Stage-IV) illustrate our core logic: leveraging high-quality camera matrices provided by SpectacularAI [1] and aligning sequence coordinates to the scene’s world frame. For detailed explanations, please refer to Sec. 3.

former moving inside the reconstructed scene, with each device providing an independent camera coordinate system. The goal is to convert these raw dual-view videos into temporally aligned and metrically consistent per-frame human and camera information, which will serve as the foundation for subsequent calibration and motion optimization.

Firstly, we use SAI to obtain per-frame calibrated cameras for each view. Let v denote the view index ($v \in \{v_1, v_2\}$), and let t index time. For each view independently, SAI provides intrinsics and extrinsics $(K_v, R_{v,t}, T_{v,t})$ for every decoded frame $I_{v,t}$ in the native coordinate system of that view. Next, we extract human-related information using several off-the-shelf models: (i) YOLO [56] for person detection and proposal pruning; (ii) ViTPose [70] for 2D human keypoints with confidence scores; (iii) SAM2 [48] for person segmentation masks; (iv) PromptDA [27] to refine dual-view depths; and (v) VIMO [65] for camera space SMPL parameters. Finally, we employ a laser pointer cue for frame-level synchronization between the two camera streams. By identifying the frame index where the laser dot disappears, we temporally align both videos and slice all associated image, depth, and parameter data accordingly. This process yields synchronized dual-view RGB-D sequences with calibrated camera trajectories and per-frame human priors, providing clean inputs for subsequent sequence calibration.

3.3. Stage III: Sequence Calibration

After obtaining the static scene reconstruction in Stage 3.1 and the dual-view camera trajectories in Stage 3.2, the next step is to align all coordinate systems into a unified world frame. At this point, we have three separate coordinate systems: one for the reconstructed scene and two for each iPhone camera trajectory estimated by SAI. Since the dual-

view coordinate systems differ from the scene coordinate system only by rigid transformations, our goal is to optimize these 2 rigid transformations to unify the dual-view coordinates into the same metric, gravity-aligned world frame. The optimization process is sensitive to the initial values; therefore, it is necessary to first obtain a good initial estimate for the rigid transformations.

Get Initial Transformation from COLMAP. We register each dual-view sequence to the sparse COLMAP model constructed in Stage 3.1 using the known intrinsics K_v and background-only SIFT features \mathcal{F}_v , extracted from images with human regions removed. Matches are established through a trained vocabulary tree [52], and images are registered against the sparse COLMAP model to obtain COLMAP camera poses $(\hat{R}_{v,t}, \hat{T}_{v,t})$ in the same metric, gravity-aligned world coordinates as the scene.

To obtain the initial rigid transformation aligning the SAI camera trajectories $T_{v,t}$ with their COLMAP counterparts $\hat{T}_{v,t}$, we solve for an offset transformation $(s^{\text{off}}, R^{\text{off}}, T^{\text{off}})$ by minimizing:

$$\min_{s^{\text{off}}, R^{\text{off}}, T^{\text{off}}} \sum_{t=1}^N \|\hat{T}_t - (s^{\text{off}} R^{\text{off}} T_t + T^{\text{off}})\|_2^2, \quad (1)$$

where N is the number of frames. After centering the trajectories, we solve this minimization problem using singular value decomposition (SVD).

For gravity alignment, R^{off} is constrained to rotations about the z -axis, ensuring proper alignment of SAI trajectories with the COLMAP coordinate system.

Calibration via Multiple Constraints. While the rigid transformations obtained in the previous step provide coarse alignment between the two camera trajectories and the reconstructed scene, this initialization alone is not sufficient to achieve accurate synchronization and metric consistency.

To further refine the calibration, we jointly optimize all alignment parameters by introducing multiple geometric and photometric constraints across views. Specifically, we optimize the per-view global offsets R_v^{off} (constrained to z -axis rotations) and T_v^{off} , using the initial alignment as the starting value. The aligned camera extrinsics are:

$$R_{v,t}^{\text{ali}} = R_v^{\text{off}} R_{v,t}, \quad T_{v,t}^{\text{ali}} = R_v^{\text{off}} T_{v,t} + T_v^{\text{off}}. \quad (2)$$

The optimization minimizes a composite loss of point tracking loss, Chamfer distance, and bundle adjustment loss to ensure spatial consistency between views and the global reconstruction.

$$\mathcal{L}_{\text{calib}} = \lambda_{\text{track}} \mathcal{L}_{\text{track}} + \sum_v \lambda_{\text{ch}} d_{\text{Chamfer}} + \sum_v \lambda_{\text{ba}} \mathcal{L}_{\text{ba},v}. \quad (3)$$

Through VGGT tracking, a subset of keyframes is selected, yielding accurate dual-view pixel tracking results in the human masks region. The tracked human surface 2D pixel coordinates $q_{v,t}^{(i)}$, along with their corresponding depth values $d_{v,t}^{(i)}$, are back-projected into the world frame:

$$Q_{v,t}^{(i)} = d_{v,t}^{(i)} R_{v,t}^{\text{ali}} K_v^{-1} \begin{bmatrix} q_{v,t}^{(i)} \\ 1 \end{bmatrix} + R_{v,t}^{\text{ali}} T_{v,t}^{\text{ali}}, \quad (4)$$

To enforce track consistency between views, the following loss is minimized:

$$\mathcal{L}_{\text{track}} = \frac{1}{\sum_{v,t} |\mathcal{Q}_{v,t}|} \sum_t \sum_i \tilde{w}_t^{(i)} \|Q_{1,t}^{(i)} - Q_{2,t}^{(i)}\|_2^2, \quad (5)$$

Where $Q_{1,t}^{(i)}$ and $Q_{2,t}^{(i)}$ are the 3D back-projected coordinates of the i -th point from view 1 and view 2, respectively. The weights $\tilde{w}_t^{(i)}$ are used to control the contribution of each point based on its tracking confidence. Here $\tilde{w}_t^{(i)} = \min(w_{1,t}^{(i)}, w_{2,t}^{(i)})$ combines the VGGT confidence scores for the same point across views. The Chamfer distance term d_{Chamfer} aligns local pointclouds \mathcal{P}_v ($v \in \{v_1, v_2\}$) with the global reconstruction \mathcal{P}_g sampled from \mathcal{M}_g in Sec. 3.1, where \mathcal{P}_v is obtained by reconstructing the scene using the method from Sec. 3.1 with humans cropped by masks. The Chamfer distance is formally defined as:

$$d_{\text{Chamfer}}(\mathcal{P}_v, \mathcal{P}_g) = \frac{1}{|\mathcal{P}_v|} \sum_{p_v \in \mathcal{P}_v} \min_{p_g \in \mathcal{P}_g} \|p_v - p_g\|_2^2 + \frac{1}{|\mathcal{P}_g|} \sum_{p_g \in \mathcal{P}_g} \min_{p_v \in \mathcal{P}_v} \|p_g - p_v\|_2^2. \quad (6)$$

Finally, $\mathcal{L}_{\text{ba},v}$ ($v \in \{v_1, v_2\}$) ensures reprojection consistency for persistent matches, where the points are obtained from COLMAP image registration:

$$\mathcal{L}_{\text{ba},v} = \frac{1}{|M_v|} \sum_{(t,j) \in M_v} \|x_{v,t,j} - \pi(K_v, R_{v,t}^{\text{ali}}, T_{v,t}^{\text{ali}}, X_j)\|_2^2. \quad (7)$$

We solve Eq. (3) using the Adam [23] optimizer with gradient clipping. For yaw-only updates, R_v^{off} is parameterized by a single z -axis angle to preserve gravity alignment.

3.4. Stage IV: Motion Optimization

After obtaining calibrated dual-view trajectories and a unified scene coordinate system in Stage 3.3, we further refine the human reconstruction results to achieve accurate and temporally consistent body motions in the world frame. At this stage, both camera poses and scene geometry are fixed, allowing us to focus on optimizing the human parameters. We first triangulate dual-view 2D keypoints into world-space 3D keypoints, which serve as reliable geometric constraints across views. Then, we optimize the SMPL parameters using these triangulated 3D keypoints to recover precise body poses and translations under the unified world coordinate system.

3D Keypoint Triangulation. To triangulate the 3D keypoints $Y_{t,j}$ from their 2D projections $\{y_{v,t,j}\}$, we estimate the 3D position by minimizing the weighted reprojection error across all views:

$$\min_{Y_{t,j}} \sum_{v=1}^V c_{v,t,j} \|y_{v,t,j} - P_v Y_{t,j}\|_2^2, \quad (8)$$

where $P_v = K_v [R_{v,t} | T_{v,t}]$ is the camera projection matrix for the v -th view. The problem can be formulated as a weighted least squares optimization. Using SVD, $Y_{t,j}$ is obtained as the right singular vector corresponding to the smallest singular value of A .

World-Space SMPLify. Start from initial shape β_0 and body pose $\theta_t^{b,0}$ in Sec. 3.2, our World Frame SMPLify [30] jointly optimizes shape $\beta \in \mathbb{R}^{10}$, per-frame pose $\theta_t = \{\theta_t^g, \theta_t^b\} \in \mathbb{R}^{72}$ and root translation $\gamma_t \in \mathbb{R}^3$ by minimizing:

$$\mathcal{L}_{\text{SMPLify}} = \mathcal{L}_{3D} + \mathcal{L}_{\text{smooth}} + \mathcal{L}_{\text{prior}} + \mathcal{L}_{\text{reproj}} \quad (9)$$

We use a two-stage optimization phase to ensure the smoothness and alignment with the original dual views, which will be detailed in Supp.Mat.

4. Evaluation

4.1. Ablation Study on Loss Functions

Ablation on dataset optimization. We conduct an ablation study on four core loss functions that significantly influence performance during data optimization, as described in main paper. These loss functions include tracking loss, Chamfer distance, reprojection loss, smoothness loss and kp3d loss. To evaluate the performance under different optimization settings, we employ four metrics. First, **IoU(Intersection over Union)** measures the overlap between the rendered SMPL mask and the SAM2 [48] mask. Second, **Reproj** evaluates the pixel error between the reprojected SMPL joints and the 2D keypoints detected by VITPose [70]. Third, **Depth** error is computed as the mean squared error (MSE) between the rendered depth from SMPL parameters and the sensor depths refined by PromptDA [27]. Finally,

Jitter is quantified using the same temporal foot skating metric as MotionVAE [28]. All metrics are averaged across all sequences and views to ensure a robust evaluation.

The \mathcal{L}_{track} effectively stitches the two views together, significantly improving the overall reconstruction performance, making it highly impactful on the final results. The \mathcal{L}_{kp3d} provides 3D joint positions of the human body, and compared to the reprojection loss, it eliminates the issue of depth ambiguity, thus playing a critical role in the overall performance.

Table 2. The performance of different optimization settings.

\mathcal{L}_{track}	$\mathcal{L}_{chamfer}$	\mathcal{L}_{reproj}	\mathcal{L}_{smooth}	\mathcal{L}_{kp3d}	IoU(%) \uparrow	Reproj \downarrow	Depth \downarrow	Jitter \downarrow
\times	\checkmark	\checkmark	\checkmark	\checkmark	54.3	44.2	2.372	0.0371
\checkmark	\times	\checkmark	\checkmark	\checkmark	<u>72.5</u>	10.9	0.081	0.0131
\checkmark	\checkmark	\times	\checkmark	\checkmark	72.3	11.1	<u>0.079</u>	0.0130
\checkmark	\checkmark	\checkmark	\times	\checkmark	72.1	<u>10.4</u>	0.087	0.0160
\checkmark	\checkmark	\checkmark	\checkmark	\times	59.3	20.4	0.609	0.0126
\checkmark	\checkmark	\checkmark	\checkmark	\checkmark	73.0	9.3	0.078	<u>0.0128</u>

4.2. Comparison with Optical Mocap

Direct comparison in optical mocap studio. To evaluate the accuracy of dual view capture system, we set up furniture in a mocap studio and use a Vicon system to capture ground truth human motion. Two photographers record dual-view videos of the actor with iPhones, while the actor performs basic motions(see Fig. 8, zoom in). We record 5 sequences of one participant with 9420 frames in total. We compare the errors against optical mocap GT of: monocular model GVHMR, our dual-view optimization, and our single-view version(v1 and v2). For the single-view version, we calibrate the actor coordinates to the scene coordinates system using COLMAP and optimize the motion with reprojection, smooth, and prior losses. The optical mocap results are fitted to SMPLX parameters by Mosh [31] and synchronized to dual-view parameters with foot contact keyframes. Results are compared in chunk sizes of 100, 500, and 1000. Our dual-view method outperforms the monocular model and single-view optimization by a large margin. As the chunk length increases, our advantage becomes increasingly evident. (see Tab. 8)

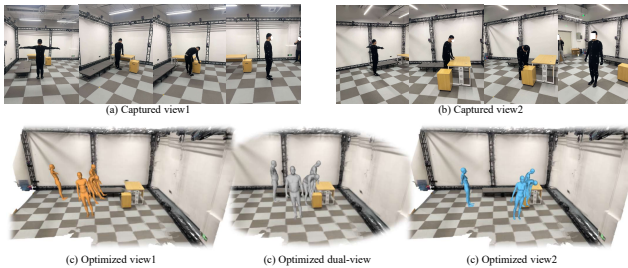


Figure 3. Optimized results in optical studio.

The advantage of dual-view over single-view lies in two key aspects: 1)dual-view effectively addresses occlusion

Table 3. Comparison among monocular model, single view optimization, with dual view optimization(ours)

Method	chunk=100		chunk=500		chunk=1000		RTE \downarrow
	WA-MPJPE \downarrow	W-MPJPE \downarrow	WA-MPJPE \downarrow	W-MPJPE \downarrow	WA-MPJPE \downarrow	W-MPJPE \downarrow	
GVHMR	66.56	123.44	124.61	333.34	179.47	593.79	1.85
Single-View V1	124.68	218.22	233.06	489.11	297.83	768.31	2.71
Single-View V2	108.31	211.83	231.41	357.22	338.42	762.80	3.65
Dual View	56.61	72.86	76.90	99.75	119.45	169.11	1.13

and self-occlusion of body joints, 2)it handles the challenging alignment of actor motion coordinates to the scene coordinates. The COLMAP estimates the camera locations for the images but suffers from depth ambiguity in the camera’s facing direction. Using a single iPhone results in large errors in the depth direction. In contrast, using two iPhones enables pixel-wise dense correspondence(see Eq. (5)), which ensures the rigid transformation between the two cameras during the optimization, and resolves the depth ambiguity in each view. **This enables a good localization of human trajectories in the scene coordinate system automatically.** Our dual view could achieve a calibration accuracy to the scene of about 5cm (human touching table in the figure), while the single view is over 30cm, measured in MeshLab by putting markers on the ground for the actor’s start and end positions.

5. Downstream Applications

In this section, we validate our capture pipeline’s effectiveness across three key applications. In Sec. 5.1, we propose a monocular human & scene reconstruction pipeline and finetune it with our captured RGBD, cameras, and SMPL annotations. In Sec. 5.2, we train several human-object interaction skills and scene-aware motion tracking with our captured motion & scene. In Sec. 5.3, we train a humanoid in simulator and deploy it to real-world robot.

5.1. Monocular Human & Scene Reconstruction

Motivation. We propose a data scheme combining RGBD data from dynamic cameras with camera and human motion parameters to train monocular human and scene reconstruction models. As no feedforward model exists, we establish a baseline using π^3 [66] for SLAM and VIMO[65] for metric-scale human motion reconstruction from monocular videos.

Implementation. To process long sequences, videos are divided into overlapping chunks, with π^3 estimating camera parameters and local point maps per chunk. Adjacent chunks are aligned using Procrustes alignment, and scale/transformations are recursively applied for global consistency. Metric scale is determined as the median ratio of SMPL to π^3 depth values. SMPL predictions are then transformed to metric world space. For details, refer to Supp. Mat. We fine-tuned two π^3 variants Tab. 4 by adding LoRA [18] layers to the camera and point decoders, supervised with the original π^3 loss. For VIMO, we froze the encoder and finetuned the decoder with MSE loss on SMPL

parameters. A human mask was used to limit supervision to the human region due to our dataset’s smaller range.

Metrics. We evaluate motion and trajectory accuracy on global coordinates using EMDB (subset 2)[22], featuring extended sequences with ground-truth trajectories and meshes. Consistent with prior work[55, 65], each sequence is split into 100-frame chunks, and 3D joint errors are measured using W-MPJPE (aligning the first two frames) and WA-MPJPE (aligning the entire segment), both in millimeters. Additionally, Root Translation Error (RTE) is reported as a percentage (%), normalized by total displacement after rigid alignment (excluding scaling).

Results. We present 3 variants in Tab. 4: the proposed baseline with the original checkpoints from π^3 [66] and VIMO [65], fine-tuning only VIMO, and fine-tuning both π^3 and VIMO. The results demonstrate that our approach significantly improves the accuracy of VIMO, as we provide paired high-quality real-world RGB sequences and ground truth SMPL parameters. Additionally, leveraging our high-quality RGB-D data and camera parameter pairs, π^3 ’s ability to predict in the world coordinate system also shows improvement. Our pipeline shows good performance on large-scale real-world videos, see Fig. 4

Table 4. Comparison of Finetuned Models on EMDB Benchmarks

Finetuned	Pi3	VIMO	EMDB		
			WA-MPJPE↓	W-MPJPE↓	RTE↓
✗	✗	✗	83.56	229.04	1.78
✗	✓	✓	82.89	222.93	1.73
✓	✓	✓	82.21	220.65	1.71

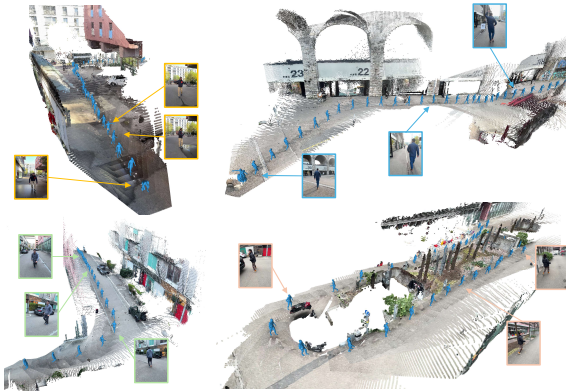


Figure 4. Quality results of proposed 4D Human & Scene Reconstruction pipeline on EMDB dataset.

5.2. Physics-based Character Animation

5.2.1. Human Object Interaction Skill Training

Motivation. We train several human-object interaction skills to demonstrate the physical realism of our approach and the scalability of our capture framework to new interaction skills. We aim to prove the efficiency and quality superiority of our framework over optical capture and monocular estimation methods.

Implementation. Following [41, 45, 64], we train physical character policies use goal-conditioned reinforcement learning to formulate character control as a Markov Decision Process (MDP) defined by states, actions, transition dynamics, a reward function r , and a discount factor γ . The reward $r_t \in \mathcal{R}$ is calculated by a style reward r_t^{style} [45] and a task reward r_t^{task} . The policies are trained to maximize the expected discounted return: $J(\pi) = \mathbb{E}_{p(\tau|\pi)} \left[\sum_{t=0}^{T-1} \gamma^t r_t \right]$, where T is the episode length, $\gamma \in [0, 1]$ is the discount factor, and r_t is the reward at time step t . We use the widely adopted Proximal Policy Optimization (PPO) algorithm [53] to train the control policy model.

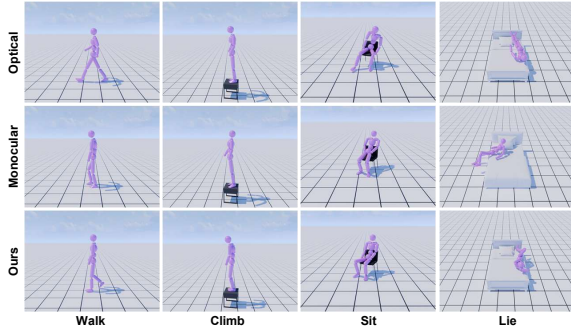
Following [13, 41, 64], we train a set of human object interaction skills in simulator [37], including *follow*, *climb*, *sit*, and *lie*. These common interaction skills are designed to guide the character’s root joint to reach specific target positions in 3D environments while maintaining physically realistic and motion diversity. We train these four common skills on 3 different input data: optical captured, which are collected from AMASS [36] and SAMP [12] following TokenHSI [41]; ours, by segmenting the reconstructed motions into skill clips; monocular, by using the motion predicted by GVHMR [54] which is commonly used in humanoid reference motion prediction[17, 67], segmented with the same temporal slices as ours. We also train 2 extra interaction skills which have not been implemented in previous physics-based human object interaction papers: Prone and Support. We will illustrate the observation, reward designs, and the training details of each skill in Supp.Mat.

Metrics. We follow [12, 68] that uses *Success Rate* and *Contact Error* as the main metrics to measure the quality of interactions quantitatively. Success Rate records the percentage of trials that humanoids successfully complete the contact within a certain threshold. We follow [13, 40, 68] in setting the thresholds for various actions: 20cm for Sit, Follow, and Climb; 30cm for Lie and Prone; and 10cm for Support. For Support, the error is defined as the distance from the object surface center to the hand center, while also taking into account the distance between the two feet. Please see details in Supp.Mat. We evaluate motion diversity using Average Pairwise Distance (APD) [7], which measures the average pairwise distance between joint rotations and positions in generated samples. Higher APD values indicate greater diversity.

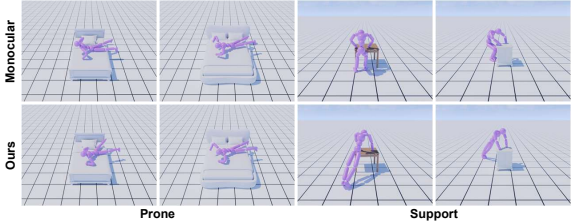
Results. We can find in Tab. 5, for skills such as Follow, Climb, and Sit, the inherent difficulty is relatively low, and all three data settings achieve good results, very close to 100%. Although the quality of our data is slightly inferior to optically captured data, we provide more variety of task completion trajectories and motion diversities, which contribute to improve task performance. To prove this, we ablate on skills trained with different data proportions. 1X and 2X indicate the ratio of the number of clips relative to the

Table 5. Comparison of data duration, Success Rate, Contact Error, and APD for different skills among 3 data settings.

Task	Data	Clips	Duration (min)	Rate (%) \uparrow	Error (cm) \downarrow	APD \uparrow
Follow	Optical Mocap	12	1.59	99.9	6.0	20.17 \pm 0.19
	Ours 1X	12	1.48	99.9	6.7	18.42 \pm 0.22
	Ours 2X	24	3.06	99.7	6.8	18.45 \pm 0.17
	Ours Full	148	22.43	99.8	6.2	19.69 \pm 0.32
	Monocular	148	22.43	98.0	7.2	19.85 \pm 0.39
Climb	Optical Mocap	7	0.28	99.9	2.7	22.03 \pm 0.30
	Ours 1X	7	0.54	99.8	1.8	22.77 \pm 0.29
	Ours 2X	14	0.97	99.9	1.8	20.72 \pm 0.30
	Ours Full	21	1.54	99.9	1.8	22.22 \pm 0.27
	Monocular	21	1.54	99.2	1.8	21.34 \pm 0.38
Sit	Optical Mocap	20	4.08	98.0	5.5	16.07 \pm 0.39
	Ours 1X	20	2.11	99.8	5.4	14.35 \pm 0.27
	Ours 2X	40	4.47	99.9	5.1	14.46 \pm 0.24
	Ours Full	80	8.05	99.9	4.7	15.90 \pm 0.51
	Monocular	80	8.05	98.4	5.7	15.80 \pm 0.51
Lie	Optical Mocap	10	2.52	89.0	17.5	8.76 \pm 0.14
	Ours 1X	10	0.99	85.3	20.2	7.43 \pm 0.10
	Ours 2X	20	2.32	86.3	19.8	8.27 \pm 0.06
	Ours Full	39	4.25	89.4	18.8	8.57 \pm 0.10
	Monocular	39	4.25	81.2	21.0	8.14 \pm 0.10
Prone	Ours Full	3	0.26	75.4	16.5	17.58 \pm 0.69
	Monocular	3	0.26	71.2	16.5	16.18 \pm 0.30
Support	Ours Full	8	0.97	66.0	4.9	21.08 \pm 0.59
	Monocular	8	0.97	20.6	6.4	20.94 \pm 0.48



(a) Qualitative comparison on 4 basic skills.



(b) Qualitative comparison on 2 additional skills.

optical capture data. On the 4 common skills, we observe a general trend where increased data amount leads to improvements in success rate, contact error, and APD metrics.

Our new implemented 2 extra skills, Prone and Support, demonstrate the versatility of our data collection pipeline. First, these new skills highlight the ability of our approach to generalize to novel interaction tasks. Second, the Support skill significantly increases the level of difficulty. Unlike other tasks, where a humanoid only needs to walk or offload the full body weight onto furniture surface, Support requires the hands to bear the weight of the body while the feet remain close together, demanding much higher accuracy in reference motion generation. This experiment shows that our approach outperforms monocular estimation methods by a large margin, particularly for high-difficulty interaction skills. The success rate trained on monocular estimated

Table 6. Quantitative evaluation of scene-aware motion tracking and dataset statistics across four 3D scenes.

Scene	Clips	Duration (min)	Status	Rate (%)	Eps. Len. (s)
a	14	12.31	Succ.	87.2	9.97 \pm 0.21
			Fail.	12.8	3.94 \pm 2.10
b	6	3.62	Succ.	96.7	9.99 \pm 0.12
			Fail.	3.3	4.16 \pm 2.38
c	12	7.87	Succ.	95.9	9.98 \pm 0.17
			Fail.	4.1	5.43 \pm 2.18
d	7	5.06	Succ.	90.4	9.96 \pm 0.21
			Fail.	9.6	4.44 \pm 1.92

motions degrades to only 20% in Tab. 5. In Fig. 5b, we can see policy trained on motion estimated from monocular models could not perform standard Support skill.

5.2.2. Scene-aware Motion Tracking

Motivation. Recent works [33, 34, 46, 57–59, 71] suggest that solving complex tasks requires pre-training on large-scale human motion data via motion tracking objectives, in order to obtain reusable and generalizable skill priors. However, existing motion tracking frameworks are mainly built for human-only [32] or single-object interaction [69] scenarios, primarily because current public datasets are concentrated in these settings. We argue that motion tracking pre-training on diverse 3D scenes is equally important, as it also provides rich priors—such as navigation, interaction, and long-horizon task execution. In this work, we mitigate this gap by: 1) proposing a scene-aware motion tracking framework, and 2) supporting it with high-fidelity paired 3D human-scene data captured by our EmbodMocap system.

Implementation. We extend MimicKit [43] by incorporating the height map into the observation space to achieve scene-aware tracking (details in the Supp. Mat.). For training, we use four 3D scenes, each containing several minutes of motion clips, and train one policy per scene to track all the motion clips in that scene.

Metrics. Policies are evaluated using a success rate metric: an episode is initialized from a random frame and run for 10s, and is considered successful if tracking exceeds 8s. For each scene, 3,072 episodes are used to compute average success, failure rates, and episode length statistics.

Results. The quantitative results in Tab. 6 demonstrate that our data is simulation-ready, enabling the training of scene-aware tracking policies with high success rates. The qualitative results, shown in Fig. 6, further illustrate that the policies not only successfully track the motions but also adapt to subtle imperfections present in the data.

5.3. Real-world Humanoid Robot Control

Motivation. Learning from human videos [2, 47, 67] has emerged as a crucial paradigm for humanoid robots to learn motor skills at scale. In this section, we demonstrate how EmbodMocap contributes to this paradigm by enabling accurate reconstruction of humans and their interacting 3D en-

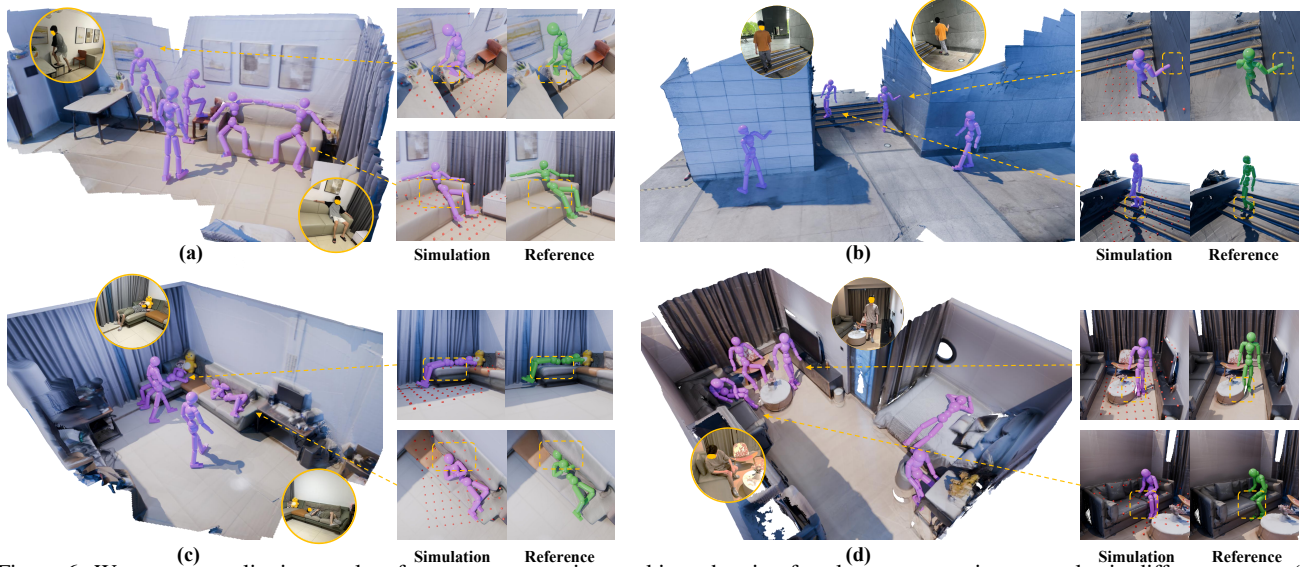


Figure 6. We present qualitative results of scene-aware motion tracking, showing four long-term motion examples in different scenes (a, b, c, and d), including daily indoor and outdoor interactions such as walking, sitting, lying, stair climbing, and touching. Our motion tracking framework not only accurately tracks the reference motion but also ensures physical realism, resolving subtle issues, such as interpenetration and floating artifacts, present in the reference data (see zoomed-in views on the right).

vironments from videos, while preserving accurate contact information.

Implementation. We capture videos of humans performing ground-contact-rich motions, including locomotion and challenging cartwheels that require precise hand-ground contact. EmbodMocap is then used for real-to-sim reconstruction. The produced motions are used to train a single tracking policy via sim-to-real RL with domain randomization using BeyondMimic [26].

Results. We deploy the policy on a real-world High Torque Hi humanoid robot with 21 joint DoF and a height of 80cm. As shown in Fig. 7, the robot successfully replicates human motions from videos, demonstrating that EmbodMocap produces data of sufficient quality for humanoid robot control.

5.4. Ablation Study on Loss Functions

Ablation on dataset optimization. We conduct an ablation study on four core loss functions that significantly influence performance during data optimization, as described in main paper. These loss functions include tracking loss, Chamfer distance, reprojection loss, smoothness loss and kp3d loss. To evaluate the performance under different optimization settings, we employ four metrics. First, **IoU(Intersection over Union)** measures the overlap between the rendered SMPL mask and the SAM2 [48] mask. Second, **Reproj** evaluates the pixel error between the reprojected SMPL joints and the 2D keypoints detected by VITPose [70]. Third, **Depth** error is computed as the mean squared error (MSE) between the rendered depth from SMPL parameters and the sensor depths refined by PromptDA [27]. Finally,

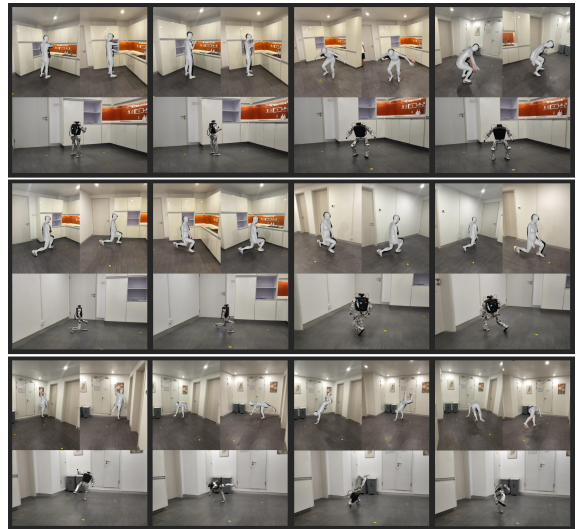


Figure 7. A real-world humanoid robot imitating human motions depicted in videos.

Jitter is quantified using the same temporal foot skating metric as MotionVAE [28]. All metrics are averaged across all sequences and views to ensure a robust evaluation.

The \mathcal{L}_{track} effectively stitches the two views together, significantly improving the overall reconstruction performance, making it highly impactful on the final results. The \mathcal{L}_{kp3d} provides 3D joint positions of the human body, and compared to the reprojection loss, it eliminates the issue of depth ambiguity, thus playing a critical role in the overall performance.

Table 7. The performance of different optimization settings.

\mathcal{L}_{track}	$\mathcal{L}_{chamfer}$	\mathcal{L}_{reproj}	\mathcal{L}_{smooth}	\mathcal{L}_{kp3d}	IoU(%) \uparrow	Reproj \downarrow	Depth \downarrow	Jitter \downarrow
\times	\checkmark	\checkmark	\checkmark	\checkmark	54.3	44.2	2.372	0.0371
\checkmark	\times	\checkmark	\checkmark	\checkmark	<u>72.5</u>	10.9	0.081	0.0131
\checkmark	\checkmark	\times	\checkmark	\checkmark	72.3	11.1	<u>0.079</u>	0.0130
\checkmark	\checkmark	\checkmark	\times	\checkmark	72.1	<u>10.4</u>	0.087	0.0160
\checkmark	\checkmark	\checkmark	\checkmark	\times	59.3	20.4	0.609	0.0126
\checkmark	\checkmark	\checkmark	\checkmark	\checkmark	73.0	9.3	0.078	<u>0.0128</u>

5.5. Comparison with Optical Mocap

Direct comparison in optical mocap studio. To evaluate the accuracy of dual view capture system, we set up furniture in a mocap studio and use a Vicon system to capture ground truth human motion. Two photographers record dual-view videos of the actor with iPhones, while the actor performs basic motions(see Fig. 8, zoom in). We record 5 sequences of one participant with 9420 frames in total. We compare the errors against optical mocap GT of: monocular model GVHMR, our dual-view optimization, and our single-view version(v1 and v2). For the single-view version, we calibrate the actor coordinates to the scene coordinates system using COLMAP and optimize the motion with reprojection, smooth, and prior losses. The optical mocap results are fitted to SMPLX parameters by Mosh [31] and synchronized to dual-view parameters with foot contact keyframes. Results are compared in chunk sizes of 100, 500, and 1000. Our dual-view method outperforms the monocular model and single-view optimization by a large margin. As the chunk length increases, our advantage becomes increasingly evident. (see Tab. 8)

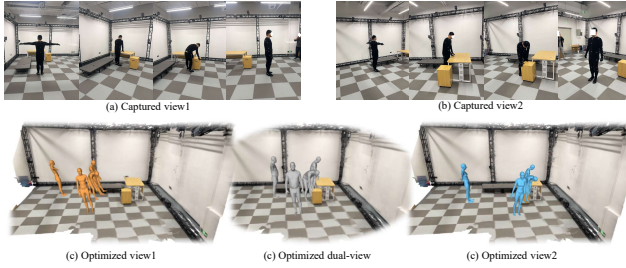


Figure 8. Optimized results in optical studio.

Table 8. Comparison among monocular model, single view optimization, with dual view optimization(ours)

Method	chunk=100		chunk=500		chunk=1000		RTE \downarrow
	WA-MPJPE \downarrow	W-MPJPE \downarrow	WA-MPJPE \downarrow	W-MPJPE \downarrow	WA-MPJPE \downarrow	W-MPJPE \downarrow	
GVHMR	66.56	123.44	124.61	333.34	179.47	593.79	1.85
Single-View V1	124.68	218.22	233.06	489.11	297.83	768.31	2.71
Single-View V2	108.31	211.83	231.41	357.22	338.42	762.80	3.65
Dual View	56.61	72.86	76.90	99.75	119.45	169.11	1.13

The advantage of dual-view over single-view lies in two key aspects: 1) dual-view effectively addresses occlusion and self-occlusion of body joints, 2) it handles the challenging alignment of actor motion coordinates to the scene coordinates. The COLMAP estimates the camera locations for the images but suffers from depth ambiguity in the

camera’s facing direction. Using a single iPhone results in large errors in the depth direction. In contrast, using two iPhones enables pixel-wise dense correspondence(see Eq. (5)), which ensures the rigid transformation between the two cameras during the optimization, and resolves the depth ambiguity in each view. ***This enables a good localization of human trajectories in the scene coordinate system automatically.*** Our dual view could achieve a calibration accuracy to the scene of about 5cm (human touching table in the figure), while the single view is over 30cm, measured in MeshLab by putting markers on the ground for the actor’s start and end positions.

6. Conclusion

We propose EmbodMocap, a portable and affordable framework for capturing high-quality 4D human & scene data using only two iPhones. Our method enables scalable, metrically accurate reconstruction of human motion and scenes mesh in diverse real-world environments. Through applications in monocular human-scene reconstruction, physics-based character animation, and humanoid robot motion control, we demonstrate the effectiveness and scalability of our approach. By lowering the barrier for embodied AI research, EmbodMocap opens new opportunities for real-world applications. We will discuss the limitations in Supp.Mat.

7. Limitations and Future Work.

Our data collection pipeline encounters limitations in specific scenarios. For example, it fails to record depth when the distance exceeds the range of the iPhone LiDAR sensor (approximately 5 meters). Additionally, it struggles with scenes dominated by moving objects, which degrade the results of the SLAM SDK [1]. Extremely bright lighting conditions can also cause COLMAP failures, leading to incorrect registration. Future work could integrate more robust structure-from-motion tools, such as H-Loc [50], to improve reliability. Moreover, incorporating automatic synchronization APPs on iPhone could further reduce human effort.

References

- [1] Spectacular ai sdk. <https://www.spectacularai.com>, 2021. 3, 4, 10, 1
- [2] Arthur Allshire, Hongsuk Choi, Junyi Zhang, David McAllister, Anthony Zhang, Chung Min Kim, Trevor Darrell, Pieter Abbeel, Jitendra Malik, and Angjoo Kanazawa. Visual imitation enables contextual humanoid control. *arXiv:2505.03729*, 2025. 3, 8
- [3] Qingwei Ben, Feiyu Jia, Jia Zeng, Juntong Dong, Dahua Lin, and Jiangmiao Pang. Homie: Humanoid loco-manipulation with isomorphic exoskeleton cockpit. *arXiv:2502.13013*, 2025. 3
- [4] Federica Bogo, Angjoo Kanazawa, Christoph Lassner, Peter Gehler, Javier Romero, and Michael J Black. Keep it smpl: Automatic estimation of 3d human pose and shape from a single image. In *ECCV*, 2016. 2
- [5] Brian Curless and Marc Levoy. A volumetric method for building complex models from range images. In *CGIT*, 1996. 3
- [6] Yudi Dai, YiTai Lin, XiPing Lin, Chenglu Wen, Lan Xu, Hongwei Yi, Siqi Shen, Yuexin Ma, and Cheng Wang. Sloper4d: A scene-aware dataset for global 4d human pose estimation in urban environments. In *CVPR*, 2023. 2, 3
- [7] Zhiyang Dou, Xuelin Chen, Qingnan Fan, Taku Komura, and Wenping Wang. C-ase: Learning conditional adversarial skill embeddings for physics-based characters. In *SIGGRAPH*, 2023. 7
- [8] Shubham Goel, Georgios Pavlakos, Jathushan Rajasegaran, Angjoo Kanazawa, and Jitendra Malik. Humans in 4d: Reconstructing and tracking humans with transformers. In *ICCV*, 2023. 2
- [9] Riza Alp Güler, Natalia Neverova, and Iasonas Kokkinos. Densepose: Dense human pose estimation in the wild. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 7297–7306, 2018. 2
- [10] Félix G. Harvey, Mike Yurick, Derek Nowrouzezahrai, and Christopher Pal. Robust motion in-betweening. 2020. 2
- [11] Mohamed Hassan, Vasileios Choutas, Dimitrios Tzionas, and Michael J Black. Resolving 3d human pose ambiguities with 3d scene constraints. In *ICCV*, 2019. 2, 3
- [12] Mohamed Hassan, Duygu Ceylan, Ruben Villegas, Jun Saito, Jimei Yang, Yi Zhou, and Michael J Black. Stochastic scene-aware motion prediction. In *ICCV*, 2021. 7
- [13] Mohamed Hassan, Yunrong Guo, Tingwu Wang, Michael Black, Sanja Fidler, and Xue Bin Peng. Synthesizing physical character-scene interactions. In *SIGGRAPH*, 2023. 7
- [14] Tairan He, Zhengyi Luo, Xialin He, Wenli Xiao, Chong Zhang, Weinan Zhang, Kris Kitani, Changliu Liu, and Guanya Shi. Omnih2o: Universal and dexterous human-to-humanoid whole-body teleoperation and learning. *arXiv:2406.08858*, 2024. 3
- [15] Tairan He, Zhengyi Luo, Wenli Xiao, Chong Zhang, Kris Kitani, Changliu Liu, and Guanya Shi. Learning human-to-humanoid real-time whole-body teleoperation. *arXiv:2403.04436*, 2024. 3
- [16] Tairan He, Wenli Xiao, Toru Lin, Zhengyi Luo, Zhenjia Xu, Zhenyu Jiang, Jan Kautz, Changliu Liu, Guanya Shi, Xiaolong Wang, Linxi Fan, and Yuke Zhu. Hover: Versatile neural whole-body controller for humanoid robots. *arXiv:2410.21229*, 2024. 3
- [17] Tairan He, Jiawei Gao, Wenli Xiao, Yuanhang Zhang, Zi Wang, Jiashun Wang, Zhengyi Luo, Guanqi He, Nikhil Sobanbabu, Chaoyi Pan, Zeji Yi, Guannan Qu, Kris Kitani, Jessica Hodgins, Linxi "Jim" Fan, Yuke Zhu, Changliu Liu, and Guanya Shi. Asap: Aligning simulation and real-world physics for learning agile humanoid whole-body skills. *arXiv preprint arXiv:2502.01143*, 2025. 3, 7
- [18] Edward J Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, Weizhu Chen, et al. Lora: Low-rank adaptation of large language models. *ICLR*, 1(2):3, 2022. 6
- [19] 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 *CVPR*, 2022. 2, 3
- [20] Mazeyu Ji, Xuanbin Peng, Fangchen Liu, Jialong Li, Ge Yang, Xuxin Cheng, and Xiaolong Wang. Ex-body2: Advanced expressive humanoid whole-body control. *arXiv:2412.13196*, 2024. 3
- [21] Angjoo Kanazawa, Michael J. Black, David W. Jacobs, and Jitendra Malik. End-to-end recovery of human shape and pose. In *CVPR*, 2018. 2
- [22] Manuel Kaufmann, Jie Song, Chen Guo, Kaiyue Shen, Tianjian Jiang, Chengcheng Tang, Juan José Zárate, and Otmar Hilliges. Emdb: The electromagnetic database of global 3d human pose and shape in the wild. In *ICCV*, 2023. 2, 3, 7
- [23] Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv:1412.6980*, 2014. 5
- [24] Muhammed Kocabas, Nikos Athanasiou, and Michael J. Black. VIBE: Video inference for human body pose and shape estimation. In *CVPR*, 2020. 2
- [25] Muhammed Kocabas, Chun-Hao P. Huang, Joachim Tesch, Lea Müller, Otmar Hilliges, and Michael J. Black. SPEC: Seeing people in the wild with an estimated camera. In *ICCV*, 2021. 2
- [26] Qiayuan Liao, Takara E Truong, Xiaoyu Huang, Guy Tevet, Koushil Sreenath, and C Karen Liu. Beyondmimic: From motion tracking to versatile humanoid control via guided diffusion. *arXiv preprint arXiv:2508.08241*, 2025. 2, 9
- [27] Haotong Lin, Sida Peng, Jingxiao Chen, Songyou Peng, Jiaming Sun, Minghuan Liu, Hujun Bao, Jiashi Feng, Xiaowei Zhou, and Bingyi Kang. Prompting depth anything for 4k resolution accurate metric depth estimation. In *Proceedings of the Computer Vision and Pattern Recognition Conference*, pages 17070–17080, 2025. 3, 4, 5, 9
- [28] Hung Yu Ling, Fabio Zinno, George Cheng, and Michiel Van De Panne. Character controllers using motion vaes. *TOG*, 2020. 6, 9
- [29] Zhizheng Liu, Joe Lin, Wayne Wu, and Bolei Zhou. Joint optimization for 4d human-scene reconstruction in the wild. *arXiv:2501.02158*, 2025. 2
- [30] Matthew Loper, Naureen Mahmood, Javier Romero, Gerard Pons-Moll, and Michael J. Black. Smpl: a skinned multi-person linear model. *TOG*, 2015. 5

- [31] Matthew M. Loper, Naureen Mahmood, and Michael J. Black. MoSh: Motion and shape capture from sparse markers. *ACM Transactions on Graphics, (Proc. SIGGRAPH Asia)*, 33(6):220:1–220:13, 2014. 6, 10
- [32] Zhengyi Luo, Jinkun Cao, Kris Kitani, Weipeng Xu, et al. Perpetual humanoid control for real-time simulated avatars. In *ICCV*, 2023. 3, 8
- [33] Zhengyi Luo, Jinkun Cao, Josh Merel, Alexander Winkler, Jing Huang, Kris Kitani, and Weipeng Xu. Universal humanoid motion representations for physics-based control. *arXiv:2310.04582*, 2023. 3, 8
- [34] Zhengyi Luo, Jinkun Cao, Sammy Christen, Alexander Winkler, Kris Kitani, and Weipeng Xu. Grasping diverse objects with simulated humanoids. *arXiv:2407.11385*, 2024. 8
- [35] Lingni Ma, Yuting Ye, Fangzhou Hong, Vladimir Guзов, Yifeng Jiang, Rowan Postyneni, Luis Pesqueira, Alexander Gamino, Vijay Baiyya, Hyo Jin Kim, et al. Nymeria: A massive collection of multimodal egocentric daily motion in the wild. In *ECCV*, 2024. 2, 3
- [36] Naureen Mahmood, Nima Ghorbani, Nikolaus F Troje, Gerard Pons-Moll, and Michael J Black. Amass: Archive of motion capture as surface shapes. In *ICCV*, 2019. 2, 7
- [37] Viktor Makoviychuk, Lukasz Wawrzyniak, Yunrong Guo, Michelle Lu, Kier Storey, Miles Macklin, David Hoeller, Nikita Rudin, Arthur Allshire, Ankur Handa, et al. Isaac gym: High performance gpu-based physics simulation for robot learning. *arXiv:2108.10470*, 2021. 7
- [38] Lea Müller, Hongsuk Choi, Anthony Zhang, Brent Yi, Jiten-dra Malik, and Angjoo Kanazawa. Reconstructing people, places, and cameras. In *Proceedings of the Computer Vision and Pattern Recognition Conference*, pages 21948–21958, 2025. 2
- [39] Riku Murai, Eric Dexheimer, and Andrew J Davison. Mast3r-slam: Real-time dense slam with 3d reconstruction priors. In *Proceedings of the Computer Vision and Pattern Recognition Conference*, pages 16695–16705, 2025. 3
- [40] Liang Pan, Jingbo Wang, Buzhen Huang, Junyu Zhang, Hao-fan Wang, Xu Tang, and Yangang Wang. Synthesizing physically plausible human motions in 3d scenes. In *3DV*, 2024. 7
- [41] Liang Pan, Zeshi Yang, Zhiyang Dou, Wenjia Wang, Buzhen Huang, Bo Dai, Taku Komura, and Jingbo Wang. Tokenhsi: Unified synthesis of physical human-scene interactions through task tokenization. In *CVPR*, 2025. 3, 7
- [42] Georgios Pavlakos, Vasileios Choutas, Nima Ghorbani, Timo Bolkart, Ahmed A. A. Osman, Dimitrios Tzionas, and Michael J. Black. Expressive body capture: 3D hands, face, and body from a single image. In *CVPR*, 2019. 2
- [43] Xue Bin Peng. Mimickit: A reinforcement learning framework for motion imitation and control, 2025. 8
- [44] Xue Bin Peng, Pieter Abbeel, Sergey Levine, and Michiel Van de Panne. Deepmimic: Example-guided deep reinforcement learning of physics-based character skills. *TOG*, 2018. 2, 3
- [45] Xue Bin Peng, Ze Ma, Pieter Abbeel, Sergey Levine, and Angjoo Kanazawa. Amp: Adversarial motion priors for stylized physics-based character control. *TOG*, 2021. 3, 7
- [46] Xue Bin Peng, Yunrong Guo, Lina Halper, Sergey Levine, and Sanja Fidler. Ase: Large-scale reusable adversarial skill embeddings for physically simulated characters. *TOG*, 2022. 3, 8
- [47] Ri-Zhao Qiu, Shiqi Yang, Xuxin Cheng, Chaitanya Chawla, Jialong Li, Tairan He, Ge Yan, David J Yoon, Ryan Hoque, Lars Paulsen, et al. Humanoid policy~ human policy. *arXiv preprint arXiv:2503.13441*, 2025. 8
- [48] Nikhila Ravi, Valentin Gabeur, Yuan-Ting Hu, Ronghang Hu, Chaitanya Ryali, Tengyu Ma, Haitham Khedr, Roman Rädle, Chloe Rolland, Laura Gustafson, et al. Sam 2: Segment anything in images and videos. *arXiv:2408.00714*, 2024. 4, 5, 9
- [49] Sara Rojas, Matthieu Armando, Bernard Ghanem, Philippe Weinzaepfel, Vincent Leroy, and Gregory Rogez. Hamst3r: Human-aware multi-view stereo 3d reconstruction. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 5027–5037, 2025. 2
- [50] Paul-Edouard Sarlin, Cesar Cadena, Roland Siegwart, and Marcin Dymczyk. From coarse to fine: Robust hierarchical localization at large scale. In *CVPR*, 2019. 10
- [51] Johannes L Schonberger and Jan-Michael Frahm. Structure-from-motion revisited. In *CVPR*, 2016. 3
- [52] Johannes Lutz Schönberger, True Price, Torsten Sattler, Jan-Michael Frahm, and Marc Pollefeys. A vote-and-verify strategy for fast spatial verification in image retrieval. In *ACCV*, 2016. 4
- [53] John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, and Oleg Klimov. Proximal policy optimization algorithms. *arXiv preprint arXiv:1707.06347*, 2017. 7
- [54] Zehong Shen, Huaijin Pi, Yan Xia, Zhi Cen, Sida Peng, Zechen Hu, Hujun Bao, Ruizhen Hu, and Xiaowei Zhou. World-grounded human motion recovery via gravity-view coordinates. In *SIGGRAPH Asia*, 2024. 2, 3, 7
- [55] Soyong Shin, Juyong Kim, Eni Halilaj, and Michael J Black. Wham: Reconstructing world-grounded humans with accurate 3d motion. In *CVPR*, 2024. 2, 7
- [56] Juan Terven, Diana-Margarita Córdova-Esparza, and Julio-Alejandro Romero-González. A comprehensive review of yolo architectures in computer vision: From yolov1 to yolov8 and yolo-nas. *Machine learning and knowledge extraction*, 2023. 4
- [57] Chen Tessler, Yunrong Guo, Ofir Nabati, Gal Chechik, and Xue Bin Peng. Maskedmimic: Unified physics-based character control through masked motion inpainting. *TOG*, 2024. 8
- [58] Chen Tessler, Yifeng Jiang, Erwin Coumans, Zhengyi Luo, Gal Chechik, and Xue Bin Peng. Maskedmanipulator: Versatile whole-body control for loco-manipulation. *arXiv preprint arXiv:2505.19086*, 2025.
- [59] Andrea Tirinzoni, Ahmed Touati, Jesse Farebrother, Mateusz Guzek, Anssi Kanervisto, Yingchen Xu, Alessandro Lazaric, and Matteo Pirodda. Zero-shot whole-body humanoid control via behavioral foundation models. In *The Thirteenth International Conference on Learning Representations*. 8
- [60] 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 *ICCV*, 2023. [2](#)

- [61] Jianyuan Wang, Minghao Chen, Nikita Karaev, Andrea Vedaldi, Christian Rupprecht, and David Novotny. Vggt: Visual geometry grounded transformer. In *CVPR*, 2025. [2](#)
- [62] Shuzhe Wang, Vincent Leroy, Yohann Cabon, Boris Chidlovskii, and Jerome Revaud. Dust3r: Geometric 3d vision made easy. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 20697–20709, 2024. [2](#)
- [63] Wenjia Wang, Yongtao Ge, Haiyi Mei, Zhongang Cai, Qingping Sun, Yanjun Wang, Chunhua Shen, Lei Yang, and Taku Komura. Zolly: Zoom focal length correctly for perspective-distorted human mesh reconstruction. In *ICCV*, 2023. [2](#)
- [64] Wenjia Wang, Liang Pan, Zhiyang Dou, Jidong Mei, Zhouyingcheng Liao, Yuke Lou, Yifan Wu, Lei Yang, Jingbo Wang, and Taku Komura. Sims: Simulating stylized human-scene interactions with retrieval-augmented script generation. *ICCV*, 2025. [3](#), [7](#)
- [65] Yufu Wang, Ziyun Wang, Lingjie Liu, and Kostas Daniilidis. Tram: Global trajectory and motion of 3d humans from in-the-wild videos. In *ECCV*. Springer, 2024. [2](#), [3](#), [4](#), [6](#), [7](#)
- [66] Yifan Wang, Jianjun Zhou, Haoyi Zhu, Wenzheng Chang, Yang Zhou, Zizun Li, Junyi Chen, Jiangmiao Pang, Chunhua Shen, and Tong He. π^3 : Scalable permutation-equivariant visual geometry learning, 2025. [6](#), [7](#)
- [67] Haoyang Weng, Yitang Li, Nikhil Sobanbabu, Zihan Wang, Zhengyi Luo, Tairan He, Deva Ramanan, and Guanya Shi. Hdmi: Learning interactive humanoid whole-body control from human videos. *arXiv:2509.16757*, 2025. [3](#), [7](#), [8](#)
- [68] 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. In *ICLR*, 2024. [7](#)
- [69] Sirui Xu, Hung Yu Ling, Yu-Xiong Wang, and Liang-Yan Gui. Intermimic: Towards universal whole-body control for physics-based human-object interactions. *arXiv:2502.20390*, 2025. [8](#)
- [70] Yufei Xu, Jing Zhang, Qiming Zhang, and Dacheng Tao. Vitpose: Simple vision transformer baselines for human pose estimation. *NeurIPS*, 2022. [4](#), [5](#), [9](#)
- [71] Heyuan Yao, Zhenhua Song, Yuyang Zhou, Tenglong Ao, Baoquan Chen, and Libin Liu. Moconvq: Unified physics-based motion control via scalable discrete representations. *ACM Transactions on Graphics (TOG)*, 43(4):1–21, 2024. [8](#)
- [72] Vickie Ye, Georgios Pavlakos, Jitendra Malik, and Angjoo Kanazawa. Decoupling human and camera motion from videos in the wild. In *CVPR*, 2023. [2](#)
- [73] Ye Yuan, Umar Iqbal, Pavlo Molchanov, Kris Kitani, and Jan Kautz. Glamr: Global occlusion-aware human mesh recovery with dynamic cameras. In *CVPR*, 2022. [2](#)
- [74] Siwei Zhang, Qianli Ma, Yan Zhang, Zhiyin Qian, Taein Kwon, Marc Pollefeys, Federica Bogo, and Siyu Tang. Ego-body: Human body shape and motion of interacting people from head-mounted devices. In *ECCV*. Springer, 2022. [2](#), [3](#)

EmbodMocap: In-the-Wild 4D Human-Scene Reconstruction for Embodied Agents

Supplementary Material

8. Human Labor Analysis

Temporal Synchronization. This step only needs the operator to identify and input the frame indices where the laser pointer’s spot disappears into a `.xlsx` file. Typically, this process takes only about 1 minute per sequence.

Skill Segmentation. Skill segmentation is only required when training physical interaction skills. The operator annotates each skill’s category, start, and end times based on the video, typically taking 0.5 to 2 minutes per sequence.

Contact Label & Optimization. In the main text, we mention that the alignment between our sequence and the scene coordinate system relies on photometric (COLMAP, pixel tracking) and geometric constraints (chamfer distance). However, this can sometimes result in alignment errors of a few centimeters, primarily due to depth inaccuracies in COLMAP’s sparse keypoints and depth errors from the iPhone sensor. To address this issue, we propose an optional post-processing solution. During data capture, we place markers in the scene and instruct the performer to begin walking from a designated marker and stop on another at the end of the sequence, standing still on the same marker. Annotating contact frame indices costs 1-2 minutes for each sequence. These markers serve as fixed reference points for alignment. In post-processing, we observe the corresponding marker positions on the reconstructed mesh and record their 3D coordinates, along with the frame indices where the performer stands on the markers. Using this information, we optimize a rigid transformation to align the center of the performer’s feet at the specified frame indices to the 3D coordinates of the markers.

Since SAI [1] could generate Z-up metric-scaled camera matrices, we define the rigid transformation in the xy-plane, defined by a rotation angle ϕ_c about the z-axis and a translation T_c . This can be represented by a homogeneous transformation matrix M :

$$M = \begin{bmatrix} R(\phi_c) & T_c \\ \mathbf{0} & 1 \end{bmatrix} = \begin{bmatrix} \cos(\phi_c) & -\sin(\phi_c) & 0 & t_x \\ \sin(\phi_c) & \cos(\phi_c) & 0 & t_y \\ 0 & 0 & 1 & t_z \\ 0 & 0 & 0 & 1 \end{bmatrix} \quad (10)$$

This matrix transform the center of lowest point on both feet to match the annotate marker. To robustly solve for the transformation parameters, we employ a gradient descent optimization, constrained by a minimizing a contact loss to

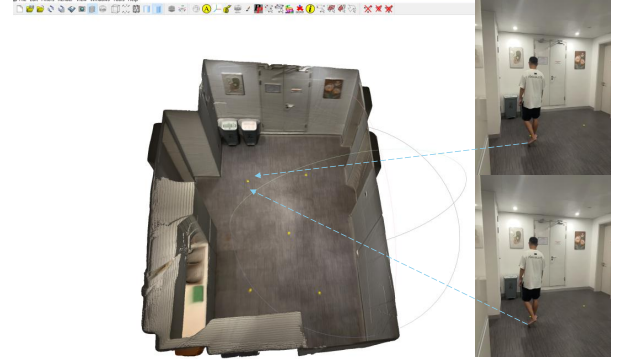


Figure 9. An example in finding the contact marker in software (e.g., Meshlab) and corresponding keyframe index(the frames selected here are just for demo).

match the contact marker:

$$\mathcal{L}_{\text{contact}} = \frac{1}{N_c} \sum_{i \in C} \left(\min_z (\mathcal{V}^{(i)}) - c_z^{(i)} \right)^2 \quad (11)$$

For SMPL parameters, the global orientation is updated as $\theta'^g = R_c \theta^g$. For translation, the pelvis’s world position is transformed as $P'_w = R_c P_w + T_c$. Re-evaluating the SMPL model with θ'^g gives the local pelvis offset P'_l , and the updated translation is $\gamma' = P'_w - P'_l$.

The updated camera rotation and translation are computed as $R'_v = R_v R_c^T$ and $T'_v = T_v - R_v R_c^T T_c$, ensuring alignment and consistency of the scene representation.

9. More Details of Monocular Human-Scene Reconstruction Pipeline

Our monocular reconstruction baseline is a modular pipeline for reconstructing 3D human pose and scene geometry from monocular video, combining two independent modules: π^3 for camera trajectory prediction and scene point cloud reconstruction, and VIMO for SMPL-based human pose estimation. To process long video sequences, π^3 divides frames into overlapping chunks, where each chunk independently predicts camera poses $T_v \in \mathbb{R}^{T \times 4 \times 4}$ and local point clouds $P_{\text{local}} \in \mathbb{R}^{T \times H \times W \times 3}$. To align these chunks into a global coordinate system, Procrustes analysis is applied to the overlapping regions of adjacent chunks. Given two point clouds $X, Y \in \mathbb{R}^{N \times 3}$, the alignment min-

imizes the error:

$$\min_{s, \mathbf{R}, \mathbf{t}} \|\mathbf{Y} - (s\mathbf{R}\mathbf{X} + \mathbf{t})\|_F^2, \quad (12)$$

where s is the scale, \mathbf{R} is the rotation matrix, and \mathbf{t} is the translation vector. Using SVD, the optimal alignment parameters are computed as:

$$\mathbf{R} = \mathbf{V}\mathbf{S}\mathbf{U}^\top, \quad s = \frac{\text{trace}(\mathbf{Y}_c^\top \mathbf{R}\mathbf{X}_c)}{\text{trace}(\mathbf{X}_c^\top \mathbf{X}_c)}, \quad \mathbf{t} = \bar{\mathbf{Y}} - s\mathbf{R}\bar{\mathbf{X}}, \quad (13)$$

where $\mathbf{X}_c, \mathbf{Y}_c$ are the centered point clouds, and \mathbf{V}, \mathbf{U} are derived from the SVD of the covariance matrix $\mathbf{H} = \mathbf{X}_c^\top \mathbf{Y}_c$. After chunk alignment, VIMO predicts SMPL parameters $(\boldsymbol{\theta}, \boldsymbol{\gamma}, \boldsymbol{\beta})$, where $\boldsymbol{\theta} \in \mathbb{R}^{T \times 72}$ represents joint rotations, $\boldsymbol{\gamma} \in \mathbb{R}^{T \times 3}$ is the root translation, and $\boldsymbol{\beta} \in \mathbb{R}^{10}$ defines body shape. Using a weak perspective camera model, SMPL vertices are projected onto the image plane as:

$$\mathbf{x}_{\text{img}} = s\mathbf{x}_v + \mathbf{t} \quad (14)$$

where s is the scaling factor proportional to $1/z$. To resolve scale ambiguity, the pipeline estimates a metric scale by matching the predicted depths of SMPL vertices z_{SMPL} (in meters) with the depths of Pi3's point cloud z_{Pi3} (in arbitrary units) on some sampled points. The scale factor is computed as:

$$s = \text{median} \left(\frac{z_{\text{Pi3}}}{z_{\text{SMPL}}} \right), \quad (15)$$

The point clouds and SMPL global orientation and translation are transformed to the world coordinate system with \mathbf{R}, \mathbf{t} following the same formula as Sec. 8.

10. More Details of Human-Object Interaction Skills

10.1. Follow Skill

Definition. The path following task requires the simulated character to move along a predefined 2D trajectory. A trajectory is represented as $\tau = \{x_{0,1}^\tau, x_{0,2}^\tau, \dots, x_{T-0,1}^\tau, x_T^\tau\}$, where $x_{0,1}^\tau$ denotes a 2D waypoint at simulation time 0.1s, and T is the episode length. For this task, T is set to 10s. The character is expected to follow the trajectory τ as accurately as possible.

Task Observation. At each simulation time step t , the character observes 10 future waypoints sampled over the next 1.0s: $\{x_t^\tau, x_{t+0.1}^\tau, \dots, x_{t+0.8}^\tau, x_{t+0.9}^\tau\}$. These waypoints are sampled at intervals of 0.1s using linear interpolation from the trajectory τ . The 2D coordinates of these waypoints form the task observation $g_t^f \in \mathbb{R}^{2 \times 10}$.

Task Reward. The reward for this task, r_t^f , is computed based on the distance between the character's current 2D

root position, $x_t^{\text{root.2d}}$, and the target waypoint, x_t^τ . The reward is defined as:

$$r_t^f = \exp(-2.0\|x_t^{\text{root.2d}} - x_t^\tau\|^2). \quad (16)$$

10.2. Sit Skill

Definition. The sitting task requires the character to position its root joint at a target 3D sitting location on an object surface. The target position is defined as 10 cm above the center of the top surface of the chair seat.

Task Observation. The observation $g_t^s \in \mathbb{R}^{38}$ includes the 3D target sitting position $\in \mathbb{R}^3$, the 3D root position $\in \mathbb{R}^3$, the root rotation $\in \mathbb{R}^6$, the 2D front-facing direction $\in \mathbb{R}^2$, and the positions of eight corner points of the object's bounding box $\in \mathbb{R}^{3 \times 8}$.

Task Reward. The sitting task reward r_t^s encourages the character to minimize the distance between its 3D root position, x_t^{root} , and the target sitting position, x_t^{tar} . It is defined as:

$$r_t^s = \begin{cases} 0.7 r_t^{\text{near}} + 0.3 r_t^{\text{far}}, & \|x_t^{\text{obj.2d}} - x_t^{\text{root.2d}}\| > 0.5, \\ 0.7 r_t^{\text{near}} + 0.3, & \text{otherwise,} \end{cases} \quad (17)$$

where r_t^{far} and r_t^{near} are defined as:

$$r_t^{\text{far}} = \exp(-2.0\|1.5 - d_t^* \cdot x_t^{\text{root.2d}}\|^2), \quad (18)$$

$$r_t^{\text{near}} = \exp(-10.0\|x_t^{\text{tar}} - x_t^{\text{root}}\|^2). \quad (19)$$

Here, $x_t^{\text{obj.2d}}$ is the 2D position of the object's root, $x_t^{\text{root.2d}}$ is the 2D linear velocity of the character's root, and d_t^* is a horizontal unit vector pointing from $x_t^{\text{root.2d}}$ to $x_t^{\text{obj.2d}}$.

10.3. Climb Skill

Definition. The climbing task requires the character to place its root joint at a target 3D climbing position on a given object. The target position is set 94 cm above the center of the top surface of the object.

Task Observation. The observation $g_t^m \in \mathbb{R}^{27}$ includes the 3D target root position $\in \mathbb{R}^3$ and the 3D coordinates of eight corner points of the object's bounding box $\in \mathbb{R}^{3 \times 8}$.

Task Reward. The climbing task reward r_t^m minimizes the 3D distance between the character's root, x_t^{root} , and the target location, x_t^{tar} . The reward is defined as:

$$r_t^m = \begin{cases} 0.5 r_t^{\text{near}} + 0.2 r_t^{\text{far}}, & \|x_t^{\text{obj.2d}} - x_t^{\text{root.2d}}\| > 0.7, \\ 0.5 r_t^{\text{near}} + 0.2 + 0.3 r_t^{\text{foot}}, & \text{otherwise,} \end{cases} \quad (20)$$

where r_t^{near} , r_t^{far} , and r_t^{foot} are defined as:

$$r_t^{\text{near}} = \exp(-10.0\|x_t^{\text{tar}} - x_t^{\text{root}}\|^2), \quad (21)$$

$$r_t^{\text{far}} = \exp(-2.0\|1.5 - d_t^* \cdot x_t^{\text{root.2d}}\|^2), \quad (22)$$

$$r_t^{\text{foot}} = \exp(-50.0\|(x_t^{\text{tar.h}} - 0.94) - x_t^{\text{foot.h}}\|^2). \quad (23)$$

Here, $x_t^{\text{tar,h}}$ is the height of the target root position, $(x_t^{\text{tar,h}} - 0.94)$ represents the height of the top surface of the target object in world coordinates, and $x_t^{\text{foot,h}}$ is the mean height of the character's feet. The reward r_t^{foot} encourages the character to lift its feet and is crucial for successful climbing.

10.4. Lie Skill

Definition. The lying task requires the character to position its root joint at a target 3D lying position on an object, typically centered on the object's surface. The character must first approach a designated standing point before transitioning into the lying position.

Task Observation. The observation $g_t^l \in \mathbb{R}^{38}$ includes the 3D target lying position $\in \mathbb{R}^3$, the 3D root position $\in \mathbb{R}^3$, the root rotation $\in \mathbb{R}^6$, the 2D front-facing direction $\in \mathbb{R}^2$, and the positions of eight corner points of the object's bounding box $\in \mathbb{R}^{3 \times 8}$. It also includes the chosen standing point $\in \mathbb{R}^3$.

Task Reward. The lying reward r_t^l combines rewards for approaching the standing point and accurately lying down:

$$r_t^l = \begin{cases} 0.6 r_t^{\text{near}} + 0.4 r_t^{\text{far}}, & \|x_t^{\text{root}} - x_t^{\text{tar}}\| > 1.5, \\ r_t^{\text{near}}, & \text{otherwise.} \end{cases} \quad (24)$$

The far reward encourages approaching the standing point:

$$r_t^{\text{far}} = 0.5 r_t^{\text{walk}} + 0.2 r_t^{\text{vel}} + 0.2 r_t^{\text{facing}} + 0.1 r_t^{\text{stand}}, \quad (25)$$

where r_t^{walk} rewards walking toward the standing point, r_t^{vel} aligns velocity, r_t^{facing} ensures proper facing direction, and r_t^{stand} rewards correct height.

The near reward focuses on lying accuracy:

$$r_t^{\text{near}} = 0.5 r_t^{\text{pos}} + 0.3 r_t^{\text{head}} + 0.2 r_t^{\text{alignment}}, \quad (26)$$

where r_t^{pos} minimizes the distance to the target, r_t^{head} aligns head height, and $r_t^{\text{alignment}}$ rewards proper body alignment.

10.5. Prone Skill

Definition. The prone task requires the character to position its root joint at a designated 3D prone position on an object, typically centered on the object's surface. Unlike the lying task, the character must face downward while maintaining alignment with the target surface.

Task Observation. The observation $g_t^p \in \mathbb{R}^{35}$ includes the 3D target prone position $\in \mathbb{R}^3$, the 3D root position $\in \mathbb{R}^3$, the root rotation $\in \mathbb{R}^6$, the 2D front-facing direction $\in \mathbb{R}^2$, and the positions of eight corner points of the object's bounding box $\in \mathbb{R}^{3 \times 8}$. These observations help guide the approach and ensure the correct orientation for prone positioning.

Task Reward. The prone reward r_t^p encourages the character to transition smoothly from moving to a prone position

while maintaining proper alignment and facing downward. The reward is defined as:

$$r_t^p = \begin{cases} 0.7 r_t^{\text{near}} + 0.3 r_t^{\text{far}}, & \|x_t^{\text{root}} - x_t^{\text{tar}}\| > 1.5, \\ r_t^{\text{near}}, & \text{otherwise.} \end{cases} \quad (27)$$

The far reward encourages approaching the target prone position:

$$r_t^{\text{far}} = 0.5 r_t^{\text{walk}} + 0.2 r_t^{\text{vel}} + 0.2 r_t^{\text{facing}} + 0.1 r_t^{\text{height}}, \quad (28)$$

where r_t^{walk} rewards moving toward the prone position, r_t^{vel} aligns velocity with the direction of motion, r_t^{facing} ensures proper facing direction, and r_t^{height} encourages maintaining an appropriate height during approach.

The near reward focuses on prone accuracy:

$$r_t^{\text{near}} = 0.6 r_t^{\text{pos}} + 0.2 r_t^{\text{alignment}} + 0.2 r_t^{\text{face_down}}, \quad (29)$$

where r_t^{pos} minimizes the distance to the prone target, $r_t^{\text{alignment}}$ ensures proper body alignment with the surface, and $r_t^{\text{face_down}}$ rewards the character for maintaining a face-down orientation.

10.6. Support Skill

Definition. The support task encourages the character to approach a target object and maintain stable interaction by placing its hands on the top surface while keeping stable foot placement and proper posture.

Task Observation. The task observation $g_t^m \in \mathbb{R}^{27}$ consists of the 3D target position of the object's top surface center ($x_t^o, z_t^o \in \mathbb{R}^3$) and the 3D coordinates of the eight corner points of the object's bounding box ($b_t \in \mathbb{R}^{3 \times 8}$).

Task Reward. The total reward r_t^m is defined as:

$$r_t^m = \begin{cases} 0.4 r_t^f + 0.6 r_t^s, & \|x_t^o - x_t^r\| > 1.5, \\ r_t^s, & \text{otherwise,} \end{cases} \quad (30)$$

$$r_t^f = 0.5 \exp(-0.5 \|x_t^o - x_t^r\|^2) \quad (31)$$

$$+ 0.5 \exp(-2.0 \|1.5 - d_t^* \cdot \dot{x}_t^r\|^2), \quad (32)$$

$$r_t^s = 0.3 r_t^h + 0.2 r_t^g + 0.15 r_t^t + 0.2 r_t^o + 0.15 r_t^z, \quad (33)$$

where r_t^f encourages the character to approach the object, and r_t^s combines five components for stable interaction:

$$r_t^h = 0.6 \exp(-20 \|z_t^h - z_t^o\|^2) \quad (34)$$

$$+ 0.4 \exp(-5 \|x_t^{h2} - x_t^o\|^2), \quad (35)$$

$$r_t^g = \exp(-50 \|z_t^f - z_g\|^2), \quad (36)$$

$$r_t^t = \exp(-10 \|x_t^{fr} - x_t^{fl}\|^2), \quad (37)$$

$$r_t^o = \exp(-2 \|1.0 - (-u_t^b)\|^2), \quad (38)$$

$$r_t^z = \exp(-10 \|z_t^r - z_t^o\|^2). \quad (39)$$

Here, x_t^o and x_t^r denote the 2D positions of the object and the character’s root, while z_t^o and z_t^r are their respective heights. x_t^{h2} and z_t^h represent the 2D position and height of the hands. Similarly, x_t^{fr} , x_t^{fl} , and z_t^f refer to the 2D positions and height of the feet, z_g is the ground height, and $-u_t^b$ is the vertical component of the body’s up direction.

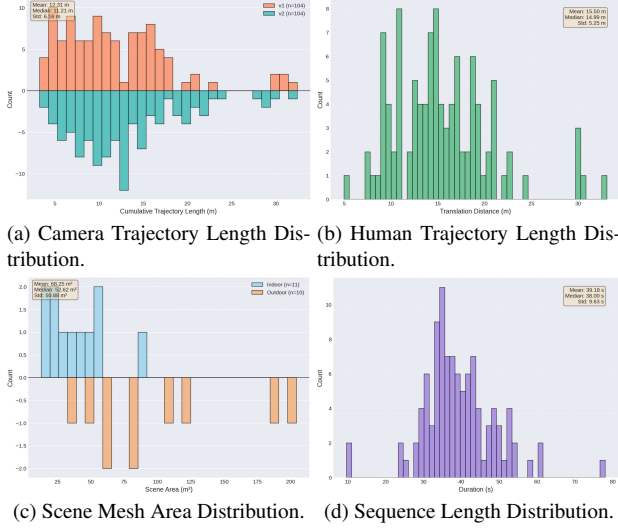


Figure 10. Statistical information of collected dataset.

Evaluation The evaluation of the Support task focuses on the agent’s ability to position its hands on the top surface of the target object and keep its feet close together. The key metric is the combined XY-plane distance and Z-axis deviation between the hands and the object’s top surface. The task is deemed successful if the hands are within predefined thresholds and the feet maintain adequate proximity for stability.

11. More Details of Scene-Aware Imitation Policy

11.1. Representations

Character Proprioception. The state s describes the proprioception of the character’s body, with features consisting of the relative positions of each link with respect to the root (designated to be the pelvis), their rotations expressed in quaternions, and their linear and angular velocities. All features are computed in the character’s local coordinate frame, with the root at the origin and the x-axis along the root link’s facing direction.

Height Map. To perceive the surrounding scene geometry, we utilize a local egocentric height map. This map is structured as an 11×11 grid spanning a $2m \times 2m$ area centered on the humanoid, resulting in a sampling interval of 0.2m. The grid is defined within the character’s local coordinate frame; consequently, the sampling points dynam-

ically translate and rotate with the humanoid’s movement and heading, consistently covering the immediate vicinity. The height values at these grid points are queried from a high-resolution underlying scene mesh (0.05m resolution) using nearest-neighbor interpolation.

Target States. The target state \hat{q} encodes the desired future motion of the character. It is constructed by sampling a short trajectory segment from the dataset spanning three consecutive future time steps: $T, T+1$, and $T+2$. For each time step, the state comprises the positions, rotations, linear velocities, and angular velocities of all body links. All features are transformed from the world frame into the simulated character’s local coordinate frame. This local frame is defined with the character’s root located at the origin and the x-axis aligned with the root link’s facing direction.

Action. Our simulated humanoid is constructed based on the SMPL body model, comprising 23 controllable joints. Each joint possesses 3 degrees of freedom (DoF), and we employ a Proportional-Derivative (PD) controller for each DoF. Consequently, the action $a \in \mathbb{R}^{69}$ generated by the policy specifies the target orientations for these PD controllers.

11.2. Reward

To encourage the character to closely reproduce the reference motion while maintaining motion naturalness, our reward function r_t is composed of two terms: a tracking reward r_t^{track} and a jitter penalty r_t^{smooth} . The tracking reward incentivizes the policy to minimize the kinematic error between the simulated character and the reference motion. The jitter penalty is introduced to suppress abnormal shaking generated when the character interacts with objects, which may be induced by instabilities in the physics simulation. The total reward is defined as:

$$r_t = r_t^{\text{track}} - r_t^{\text{smooth}}. \quad (40)$$

The tracking reward r_t^{track} is computed as the weighted sum of exponential differences across all humanoid links:

$$\begin{aligned} r_t^{\text{track}} = & w_{jp} \exp(-100\|\hat{p}_t - p_t\|^2) \\ & + w_{jr} \exp(-10\|\hat{q}_t \ominus q_t\|^2) \\ & + w_{jv} \exp(-0.1\|\hat{v}_t - v_t\|^2) \\ & + w_{j\omega} \exp(-0.1\|\hat{\omega}_t - \omega_t\|^2), \end{aligned} \quad (41)$$

where the equation penalizes the differences in translation p , rotation q , linear velocity v , and angular velocity ω for all rigid body links of the humanoid between the simulation and the reference. The jitter penalty penalizes the magnitude of the difference between consecutive actions, defined as:

$$r_t^{\text{smooth}} = \|a_t - a_{t-1}\|^2, \quad (42)$$

where a_t and a_{t-1} denote the action at the current and previous time steps, respectively. By minimizing the rate of



Figure 11. Rendered SMPL and depth images of the captured dataset in camera space.

change of the actions, the policy is incentivized to generate continuous and stable control trajectories, thereby reducing jittery behaviors.

12. More Details of Captured Dataset Used in Main Paper

We collected data from 23 scenes, each with a high-precision mesh, 104 sequences, and approximately 200,000 video frames. Each frame is accompanied by corresponding depth maps, segmentation masks, camera trajectories, and human parameters (bounding boxes, 2D keypoints, SMPL parameters).

In Fig. 10a, we present the distribution of camera trajectory lengths, which range from 4 meters to over 30 meters. In Fig. 10b, the human trajectory length distribution is shown, with performers moving between 5 meters and over 30 meters. Figure 10c illustrates the scene mesh area dis-

tribution. Indoor scenes are relatively smaller, ranging from 20 to 90 square meters, while outdoor scenes can be as large as 200 square meters. Finally, in Fig. 10d, we show the sequence length distribution, where most sequences have durations ranging from 30 to 60 seconds.

12.1. Qualitative Demonstrations

We show camera space results in Sec. 10.6 and world space results in Sec. 12.1

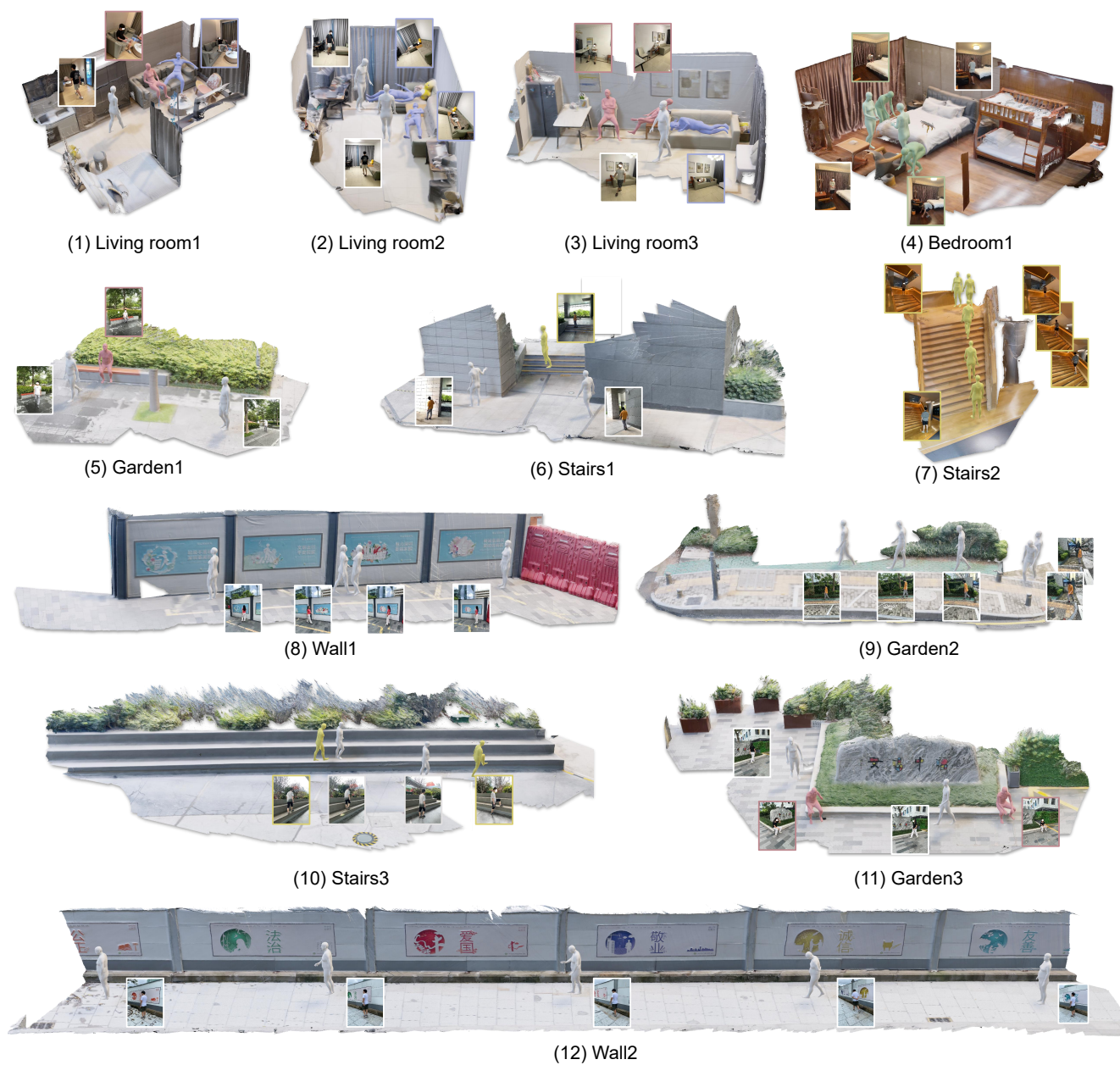


Figure 12. 3D demo of the captured dataset.