# MeshMimic: Geometry-Aware Humanoid Motion Learning through 3D Scene Reconstruction

**Qiang Zhang**[1,2,*†], **Jiahao Ma**[1,7,*], **Peiran Liu**[1,2,*], **Shuai Shi**[1,*], **Zeran Su**[1], **Zifan Wang**[2], **Jingkai Sun**[1,3], **Wei Cui**[1], **Jialin Yu**[1], **Gang Han**[1], **Wen Zhao**[1], **Pihai Sun**[1], **Kangning Yin**[6], **Jiaxu Wang**[5], **Jiahang Cao**[3], **Lingfeng Zhang**[4], **Hao Cheng**[2], **Xiaoshuai Hao**[4], **Yiding Ji**[2], **Junwei Liang**[2], **Jian Tang**[1], **Renjing Xu**[2], **Yijie Guo**[1]

[1]X-Humanoid
[2]The Hong Kong University of Science and Technology (Guangzhou)   [3]The University of Hong Kong   [4]Tsinghua University
[5]The Chinese University of Hong Kong   [6]Shanghai Jiao Tong University   [7]The Australian National University
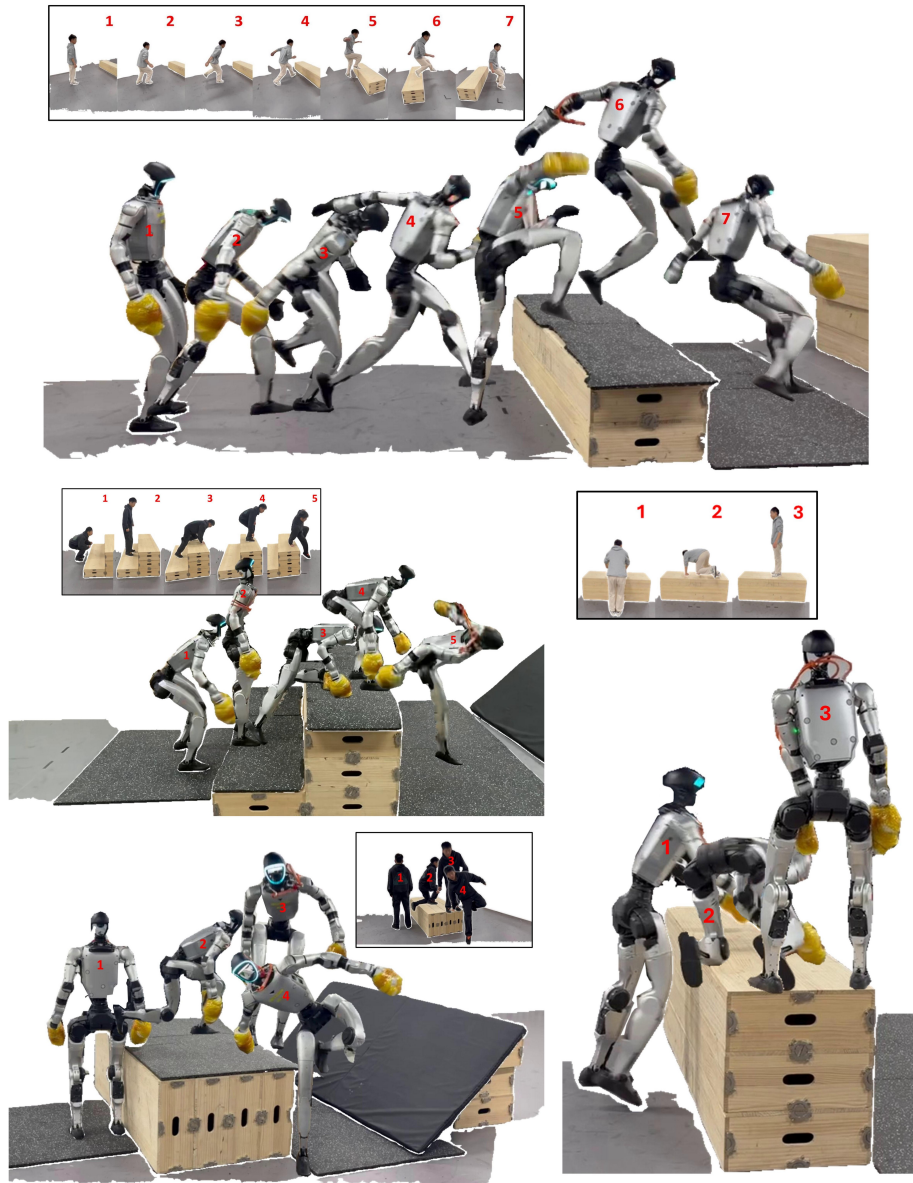[*]Equal Contribution   [†]Corresponding Author

Figure 1: **MeshMimic: monocular video-to-humanoid robots.** From ordinary *consumer monocular* videos (*no MoCap*), we reconstruct and optimize human motion together with scene geometry and contacts, then retarget the coupled motion–terrain interactions to humanoid robots that perform dynamic terrain-aware skills.

## Abstract

**Humanoid motion control has witnessed significant breakthroughs in recent years, with deep reinforcement learning (RL) emerging as a primary catalyst for achieving complex, human-like behaviors. However, the high dimensionality and intricate dynamics of humanoid robots make manual motion design impractical, leading to a heavy reliance on expensive motion capture (MoCap) data. These datasets are not only costly to acquire but also frequently lack the necessary geometric context of the surrounding physical environment. Consequently, existing motion synthesis frameworks often suffer from a decoupling of motion and scene, resulting in physical inconsistencies such as contact slippage or mesh penetration during terrain-aware tasks.**

**In this work, we present MeshMimic, an innovative framework that bridges 3D scene reconstruction and embodied intelligence to enable humanoid robots to learn coupled "motion-terrain" interactions directly from video. By leveraging state-of-the-art 3D vision models, our framework precisely segments and reconstructs both human trajectories and the underlying 3D geometry of terrains and objects. We introduce an optimization algorithm based on kinematic consistency to extract high-quality motion data from noisy visual reconstructions, alongside a contact-invariant retargeting method that transfers human-environment interaction features to the humanoid agent. Experimental results demonstrate that MeshMimic achieves robust, highly dynamic performance across diverse and challenging terrains. Our approach proves that a low-cost pipeline utilizing only consumer-grade monocular sensors can facilitate the training of complex physical interactions, offering a scalable path toward the autonomous evolution of humanoid robots in unstructured environments.**

## 1. Introduction

Humanoid robots represent one of the most formidable frontiers in robotics. Developing controllers that can manage their high-dimensional degrees of freedom (DoFs) while executing human-like movements remains a profound challenge. With the recent advancements in artificial intelligence, Reinforcement Learning (RL) has emerged as the dominant paradigm for humanoid control, enabling agents to learn complex motor skills through exploratory interaction. However, the sheer complexity of humanoid dynamics makes manual reward engineering and motion design increasingly impractical. Consequently, the field has pivoted toward motion imitation, where robots learn to replicate human behaviors using reference data.

Traditionally, humanoid imitation learning has relied heavily on Motion Capture (MoCap) data Mahmood et al. (2019). While high-fidelity, MoCap data is prohibitively expensive to acquire and restricted to controlled laboratory settings. More critically, traditional MoCap often fails to capture the **geometric context** of the surrounding environment. This decoupling of motion from the physical scene leads to significant physical inconsistencies—such as "foot skating," contact misalignment, or mesh penetration—when the robot is tasked with navigating complex, non-flat terrains. Furthermore, in many real-world scenarios, deploying inertial or optical MoCap systems is logistically impossible. Recent works like *VideoMimic* Allshire et al. (2025) have attempted to bypass MoCap by using video data; however, they often rely on coarse scene modeling and lack fine-grained contact optimization. Others, such as *OmniRetarget* Yang et al. (2025), introduce "Interaction Meshes" to refine object manipulation but are limited to simple geometric primitives and fail to generalize to irregular, large-scale terrains.
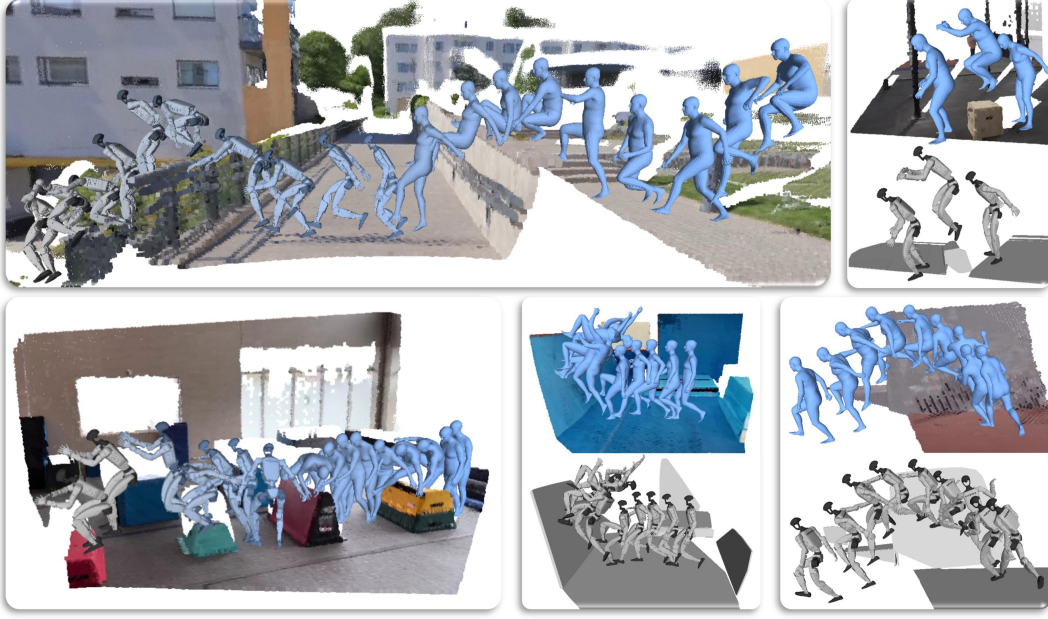
Figure 2: **MeshMimic Real-to-Sim.** In-the-wild monocular videos yield *long-horizon motions* over *complex terrains* for contact-consistent *motion–terrain* interaction learning.

The rapid evolution of computer vision offers a transformative opportunity to bridge this gap. Modern spatial representations, such as 3D Gaussian Splatting (3DGS) Kerbl et al. (2023) and Neural Radiance Fields (NeRF) Mildenhall et al. (2021), have drastically improved the quality of 3D scene reconstruction. Concurrently, foundation models like SAM3D Team et al. (2025) have enhanced the robustness of segmentation and 3D generation. Building on these developments, we posit that modern perceptual models can be harnessed to form a unified framework that recovers both high-quality human motion and the associated 3D environmental geometry directly from raw video, enabling downstream RL training.

In this paper, we propose **MeshMimic**, a novel framework that enables humanoid robots to learn complex, terrain-aware tasks directly from monocular video as shown in Fig. 1. Our system utilizes 3D vision foundation models to decouple and reconstruct human trajectories and environmental meshes from unposed video sequences. Reconstruction demonstrations are shown in Fig. 2. To handle the inherent noise in visual reconstruction, we introduce a **Kinematic Consistency Optimization** algorithm to ensure physically plausible reference motions. Furthermore, we present **MeshRetarget**, a contact-aware retargeting method that maps the intricate interactions between humans and 3D surfaces—such as stepping on uneven rocks or navigating obstacles—onto the humanoid morphology. By integrating these "physically-grounded" references into an RL pipeline, our robot learns to perceive and interact with its environment in a unified manner.

The primary contributions of this work are summarized as follows:

- **A Terrain-Aware Humanoid Motion Mimic Framework:** We present MeshMimic, an integrated framework that enables humanoid robots to learn diverse motor skills directly from monocular videos. Crucially, we introduce a **Kinematic Consistency Optimization** strategy during the motion extraction phase, which refines noisy visual pose estimations into physically plausible reference trajectories suitable for control.
- **MeshRetarget Mechanism:** We propose a novel retargeting method that explicitly addresses the morphological gap between human subjects and robots. By prioritizing geometric interaction features, MeshRetarget ensures that motions from humans of varying heights are effectively mapped to robots of different dimensions while preserving essential contact constraints.
- **Experimental Validation:** We validate our framework across a variety of highly dynamic tasks on irregular terrains. Our experiments demonstrate that MeshMimic achieves superior robustness and physical realism

compared to existing scene-agnostic baselines.

## 2. Preliminaries and Related Works

### 2.1. 3D Spatial Modeling and Environment Reconstruction

The advancement of computer vision has transitioned from sparse point-cloud representations to high-fidelity, dense geometric reconstructions. Early milestones in Structure-from-Motion (SfM) and Multi-View Stereo (MVS), exemplified by frameworks like COLMAP Fisher et al. (2021) and Schonberger and Frahm (2016), laid the foundation for spatial mapping but often struggled with textureless surfaces and dynamic occlusions common in real-world robotic environments. The field was further revolutionized by Neural Radiance Fields (NeRF) Mildenhall et al. (2021) and its accelerated variants like Instant-NGP Müller et al. (2022), which introduced differentiable volumetric representations. More recently, 3D Gaussian Splatting (3DGS) Kerbl et al. (2023) has emerged as a state-of-the-art representation, offering explicit, primitive-based modeling with real-time rendering capabilities.

However, for humanoid robotics, raw geometric reconstruction is insufficient without semantic or instance-level decomposition to distinguish navigable terrain from dynamic agents. While works like LERF Kerr et al. (2023) and ConceptFusion Jatavallabhula et al. (2023) attempted to ground foundation models into 3D spaces for robotic manipulation, they often lack the fine-grained geometric precision required for high-dynamic locomotion. The emergence of 3D-aware foundation models, such as SAM3D Team et al. (2025) and Segment Anything in High Quality (HQ-SAM) Ke et al. (2023), has enabled robust instance-level segmentation within 3D space. These models allow for the isolation of human actors from their environmental context with unprecedented accuracy. Unlike previous motion imitation works that treat the environment as a simplified static background or a flat plane, our approach leverages these 3D segmentation capabilities to extract the local geometry of the terrain. By converting these segmented instances into high-resolution collision meshes, we provide the necessary physical constraints and exteroceptive observations for downstream reinforcement learning.

### 2.2. Humanoid Motion Retargeting

Motion retargeting is the foundational process of mapping human kinematic trajectories onto a robot's morphology while preserving the semantic and physical intent of the motion. This task is inherently ill-posed due to significant discrepancies in degrees of freedom (DoFs), joint limits, and mass distributions. Early optimization-based approaches, such as GMR Araujo et al. (2025), focused on preserving geometric relationships and manifold structures to maintain motion fidelity. However, these methods often struggle with the dynamic stability required for high-dimensional humanoid control.

With the rise of large-scale data-driven methods, the field has shifted toward learning-based retargeting and control. Frameworks like PHC Luo et al. (2023) have demonstrated the efficacy of learning robust control policies from large-scale human motion data in simulation. Building upon this, OmniH2O He et al. (2024) introduced a universal system for full-body humanoid-to-humanoid and human-to-humanoid mapping, enabling real-time teleoperation and diverse skill acquisition. More recently, Spider Pan et al. (2025) pushed the boundaries of agile motion retargeting by optimizing for highly dynamic and versatile humanoid behaviors.

Despite these advancements, a critical gap remains in environment-aware retargeting. While OmniRetarget Yang et al. (2025) introduced interaction meshes to refine contacts between the robot and manipulated objects, most current frameworks—including PHC Luo et al. (2023) and OmniH2O Yang et al. (2025) — primarily focus on the agent's internal state or assume a simplified flat-ground plane. This lack of terrain awareness leads to physical inconsistencies, such as "foot skating" or penetration, when the robot interacts with non-planar geometries. Our proposed *MeshRetarget* addresses this by explicitly incorporating high-resolution reconstructed meshes of the terrain into the retargeting loop. By prioritizing contact-invariance on irregular surfaces, we ensure that the retargeted motion is not only kinematically feasible but also geometrically grounded in the

actual physical environment.

## 2.3. Humanoid Motion Tracking and Whole-Body Control

Humanoid motion tracking aims to bridge the gap between kinematic reference trajectories and dynamic execution within a physics-based simulator. While early character animation works like DeepMimic Peng et al. (2018) and ASE Peng et al. (2022) demonstrated the potential of reinforcement learning (RL) for motion imitation, humanoid robotics requires a higher degree of physical robustness and whole-body coordination (WBC). Recent advancements have shifted toward unified whole-body controllers that can handle the high-dimensional, non-linear dynamics of robotic hardware.

A significant milestone in this direction is ExBody Cheng et al. (2024) and its successor ExBody2 Ji et al. (2024), which facilitate expressive whole-body control by learning from human motion data. These frameworks emphasize the importance of capturing subtle upper-body gestures alongside stable locomotion, providing a more comprehensive imitation of human behavior than traditional gait-focused controllers. Similarly, frameworks like OmniH2O He et al. (2024) have pioneered the "Human-to-Humanoid" pipeline, enabling robots to track diverse human motions in real-time. To scale these capabilities, **BeyondMimic** Liao et al. (2025) and **VideoMimic** Allshire et al. (2025) have explored utilizing large-scale datasets and raw video inputs. Similarly, frameworks such as **Sonic** Luo et al. (2025), **kungfubot** Xie et al. (2025) Han et al. (2025) and **UniTracker** Yin et al. (2025) have demonstrated that training on massive motion libraries significantly enhances the generalization of the whole-body controller, allowing for more versatile and robust robotic behaviors across diverse scenarios.

Despite these breakthroughs, a critical limitation persists: most existing humanoid trackers are essentially "scene-agnostic." They typically operate under the assumption of a uniform, flat ground plane and lack exteroceptive awareness of the specific terrain geometry. This prevents the whole-body controller from proactively adjusting its gait or contact points when navigating obstacles—such as stepping over debris or traversing uneven slopes—that were inherent to the original human motion. *MeshMimic* addresses this challenge by integrating high-fidelity 3D scene reconstruction directly into the motion-tracking loop. By grounding the whole-body controller in the reconstructed geometry of the environment, we enable the robot to perform terrain-aware motion imitation that is physically consistent with both the human reference and the environmental constraints.

## 3. MeshMimic

The core philosophy of MeshMimic is to bridge the gap between unstructured visual observations and robust humanoid control through a comprehensive **Real-to-Sim-to-Real** pipeline (Fig. 3). This unified framework enables the robot to not only learn from human demonstrations in diverse environments but also to deploy the resulting intelligence back into the physical world.

The **Real-to-Sim** process constitutes the foundation of our data generation. Starting from a casually captured monocular video, we utilize a 3D-aware perception module to decouple the human actor from the environment. Through a joint optimization of the human trajectory and the reconstructed scene mesh, we recover metrically consistent human-scene interactions (Sec. 3.2). These refined trajectories are subsequently mapped to a humanoid morphology via our *MeshRetargeting* algorithm, which prioritizes contact-invariance and collision avoidance within the reconstructed 3D terrain (Sec. 3.3). This results in a high-fidelity, physically-grounded simulation environment where the humanoid agent can learn complex motor skills via Reinforcement Learning.

The **Sim-to-Real** phase enables the deployment of learned policies onto physical humanoid robot. By training with terrain-optimized reference trajectories, the whole-body controller (WBC) internalizes the interaction priors necessary for complex environments. This approach ensures that the robot maintains high contact fidelity and physical stability on specific non-flat surfaces, effectively replicating the intricate human-environment interactions captured in the original video.(Sec. 4.2)
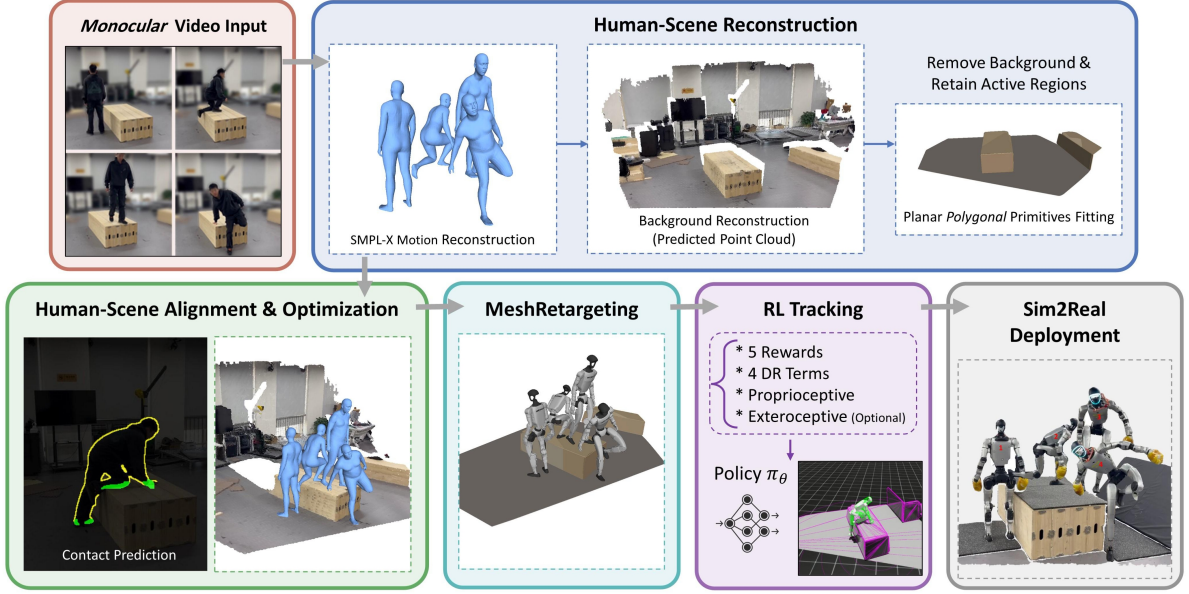
Figure 3: **MeshMimic Real-Sim-Real Pipeline.** Starting from a monocular video, we reconstruct the scene geometry and human motion, jointly align them to recover metrically consistent human–scene interactions, and retarget the refined motion to a humanoid in simulation for RL policy learning. Finally, we deploy the learned policy to the real robot enabling stable execution over challenging terrain.

## 3.1. Preprocessing

We preprocess a monocular RGB video using off-the-shelf scene reconstruction, detection/tracking, and monocular human reconstruction modules. For the environment, we run $\pi^3$ Wang et al. (2025) to reconstruct the scene and obtain per-frame depth maps $D^t$, camera poses $[R^t \mid \mathbf{t}^t]$, and a shared camera intrinsics $K$. In our scene processing, we depart from approaches that convert point clouds into dense meshes (e.g., VideoMimic) or represent the environment with simple planar primitives (e.g., CRISP Wang et al. (2025)). Instead, we approximate the scene using planar polygonal primitives. This representation effectively suppresses noisy points in dynamic reconstructions, provides a simple yet faithful scene description, and captures richer geometric structure than conventional planar primitives. For the human, we detect the target person using ViTDet Li et al. (2022) and associate the identity across frames via SAM2 Ravi et al. (2024). Given the tracked person instances, we reconstruct per-frame human body geometry and motion using SAM3D Team et al. (2025). Specifically, we follow the official SAM3D-Body pipeline: we convert the intermediate MHR representation to SMPL-X Pavlakos et al. (2019), yielding per-frame local pose parameters $\boldsymbol{\theta}^t$, body shape $\boldsymbol{\beta}$, and 3D SMPL joints $\mathbf{J}_{3D}^t \in \mathbb{R}^{J \times 3}$, together with the estimated orientation $\phi^t$ translation $\mathbf{t}^t$ in camera coordinate.

Importantly, both the scene reconstruction (camera parameters and geometry from $\pi^3$) and the monocular human reconstruction are not metrically scaled. Moreover, the recovered SMPL-X motion is expressed in the camera coordinate system and thus is not directly comparable across frames in a common world frame. In the subsequent stage, we therefore jointly optimize the human motion and scene geometry to recover a metrically consistent, world-aligned human trajectory and environment geometry.

## 3.2. Human-Scene Reconstruction

Unlike VideoMimic, which focuses on curated videos with clean motion, our Internet-crawled clips often exhibit rapid camera/subject motion, leading to jitter, blur, and occlusions. In this regime, learning-based contact prediction like BSTRO Huang et al. (2022) becomes unstable and degrades human–scene optimization. We address this with depth-edge–guided contact prediction, metric-scale human–scene alignment and joint human-scene optimization for robust contacts and metrically consistent reconstruction.
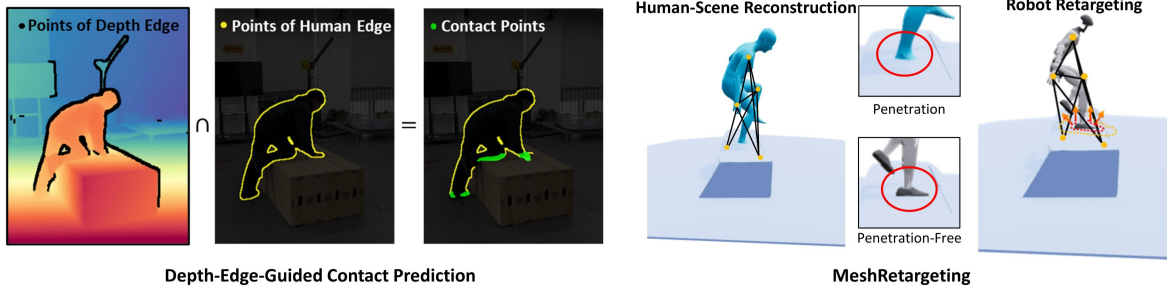
Figure 4: **Left:** Depth-edge–guided contact prediction. **Right:** MeshRetargeting optimization for penetration correction and contact-consistent retargeting.

**Depth-edge–guided contact prediction**. Visual cues such as depth edges $E_\text{depth}$ from monocular depth and silhouette boundaries $E_\text{human}$ from human segmentation are relatively stable. We exploit them to extract reliable human–scene contacts. We compute $E_\text{human}$ by applying a morphological gradient to the binary human mask, and dilate $E_\text{depth}$ to form an exclusion region $\tilde{E}_\text{depth}$ around depth discontinuities. We then define the contact band as the silhouette-boundary pixels not covered by $\tilde{E}_\text{depth}$, i.e., $\mathcal{P}_c = \{\, p \in \mathcal{P}_\text{human} \mid E_\text{human}(p) = 1 \land \tilde{E}_\text{depth}(p) = 0 \,\}$, where $p = (u, v)$ denotes a pixel, $\mathcal{P}_\text{human}$ is the set of silhouette-boundary pixels, and $E_\text{human}(p), \tilde{E}_\text{depth}(p) \in \{0, 1\}$ indicate whether $p$ lies on the human boundary or within the dilated depth-edge exclusion region, respectively. We further dilate $\mathcal{P}_c$ with a small kernel to improve robustness to projection noise. In Figure. 5.A, human-scene contacts are consistently highlighted in green across the video sequence. Finally, background points whose projections fall within this band are selected as candidate scene contacts.

**Metric-scale human–scene alignment**. We jointly optimize the human trajectory and the scene scale. Since SAM-3D Body provides a strong initialization, we keep the SMPL-X pose and shape fixed and optimize only global alignment variables: the per-frame translations $\mathbf{t}^{0:T}$ and a single scene scale $\alpha$ applied to the reconstructed point cloud. Using the metric height prior of SMPL-X as a reference, $\alpha$ corrects the scale mismatch between scene reconstruction and human-derived metric scale. Let $\mathbf{H}$ denote the SMPL-X function that outputs vertices $\mathcal{V}_h^t = \mathbf{H}(\boldsymbol{\beta}, \boldsymbol{\theta}^t, \mathbf{t}^t, \boldsymbol{\phi}^t)$. Following Yalandur Muralidhar et al. (2025), we only optimize the global translation $\mathbf{t}^t$ using an alignment loss $L_\text{align} = L_{J2d} + L_d$. Here, $L_{J2d}$ measures the 2D joint reprojection error between SMPL-X–regressed joints (projected with the camera intrinsics) and SAM3D-body 2D keypoints prediction, while $L_d$ is a symmetric Chamfer distance between camera-facing SMPL-X vertices and the metric-scale human point set. To avoid matching backside vertices to the human points, we compute $L_d$ using only camera facing vertices, selected by the angle between each vertex normal and viewing direction to be less than $65°$.

**Kinematic consistency optimization.** Although the human vertices $\mathcal{V}_h^t$ and metric-scale scene points provide a strong initialization for human–scene alignment, occlusions, inaccurate depth estimates, and noisy 2D keypoint predictions can still cause interpenetration or hovering artifacts, as well as unstable and drifting trajectories. To improve physical plausibility, we further enforce human–scene kinematic consistency by introducing additional constraints. Specifically, we jointly optimize the per-frame translations $\mathbf{t}^{0:T}$ and a single global scene scale using an alignment term together with contact, penetration, trajectory smoothness, and foot-snapping regularization:

$$L_\text{total} = \lambda_{align} L_{align} + \lambda_c L_c + \lambda_p L_p + \lambda_{sm} L_{sm} + \lambda_\text{fs} L_\text{fs}. \tag{1}$$

We detail the contact loss ($L_c$), penetration loss ($L_p$), trajectory smoothness loss ($L_{sm}$), and foot-snapping loss ($L_{fs}$) in the following.

*Contact loss ($L_c$).* We encourage predicted contacting human vertices to coincide with the estimated scene contact locations. For each frame $t$, we obtain a set of scene contact points $\{\mathbf{c}_j^t\}$ and the corresponding human vertex indices $\{i_j^t\}$. Since the reconstructed scene is re-scaled by a single global factor $\alpha$, we enforce contact

consistency by matching each scaled scene contact point $\alpha\,\mathbf{c}_j^t$ to its corresponding human vertex $\mathbf{v}_{h,i_j^t}^t \in \mathcal{V}_h^t$:

$$L_c = \frac{1}{\sum_t |\mathcal{C}^t|} \sum_t \sum_{j \in \mathcal{C}^t} \left\| \alpha\,\mathbf{c}_j^t - \mathbf{v}_{h,i_j^t}^t \right\|_2^2, \tag{2}$$

which anchors contact vertices to the scene and reduces hovering at predicted support regions.

*Penetration loss ($L_p$).* To discourage the human mesh from intersecting scene geometry, we construct a TSDF volume Curless and Levoy (1996); Newcombe et al. (2011) from the background point cloud and oriented normals. We then query the TSDF at all world-space human vertices using trilinear sampling. Denote the sampled signed distance at vertex $\mathbf{v}$ by $d(\mathbf{v})$, where $d > 0$ is outside and $d < 0$ indicates penetration. We introduce a slack $\tau$ and only penalize penetration deeper than $-\tau$, and denote penalty as $p(\mathbf{v}) = \max\big(0,\ -(d(\mathbf{v}) + \tau)\big)$. Finally, we apply a Huber-style robust penalty to stabilize gradients:

$$L_p \;=\; \frac{1}{|\mathcal{V}|} \sum_{\mathbf{v} \in \mathcal{V}} \mathrm{Huber}(p(\mathbf{v}))\,. \tag{3}$$

so shallow violations are softly corrected while deeper intersections are strongly penalized.

*Trajectory smoothness loss ($L_{sm}$).* To mitigate frame-to-frame jitter and drift in the recovered global motion, we regularize the per-frame global translation trajectory. Let $\mathbf{T}^t = \mathbf{t}_{\mathrm{cam}}^t + \mathbf{t}^t$ denote the global translation at frame $t$, and let $N$ be the number of frames in the sequence ($t = 0, \dots, N-1$). Inspired by Wang et al. (2025), we penalize both velocity and acceleration using finite differences scaled by the frame rate $f$:

$$L_{sm} \;=\; \frac{1}{N-1} \sum_{t=0}^{N-2} \left\| (\mathbf{T}^{t+1} - \mathbf{T}^t)\,f \right\|_2^2 \;+\; \frac{1}{N-2} \sum_{t=0}^{N-3} \left\| (\mathbf{T}^{t+2} - 2\mathbf{T}^{t+1} + \mathbf{T}^t)\,f \right\|_2, \tag{4}$$

which encourages temporally coherent motion while still allowing genuine fast movements.

*Foot-snapping loss ($L_{fs}$).* To reduce near-ground foot hovering (i.e., "foot snap" artifacts), we explicitly encourage foot joints to lie on the scene surface when they are already close to it. We extract foot joint positions $\mathbf{q}_f^t$ from the optimized 3D keypoints, evaluate their TSDF values $d(\mathbf{q}_f^t)$, and activate the term only within a narrow near-surface band $\mathbb{I}\big(0 < d(\mathbf{q}_f^t) \leq \tau_{\mathrm{contact}}\big)$, where $\tau_{\mathrm{contact}}$ is a contact threshold. Within this band, we penalize squared distances to pull the feet onto the surface:

$$L_{fs} \;=\; \frac{1}{N} \sum_{t,f} \mathbb{I}\big(0 < d(\mathbf{q}_f^t) \leq \tau_{\mathrm{contact}}\big)\, d(\mathbf{q}_f^t)^2. \tag{5}$$

These two terms $L_{fs}$ and $L_p$ are complementary: $L_{fs}$ reduces near-surface hovering by penalizing small *positive* TSDF distances, whereas $L_p$ discourages interpenetration by penalizing *negative* TSDF distances.

## 3.3. MeshRetargeting

Following OminiRetarget Yang et al. (2025), we preserve the spatial relationships among robot parts, manipulated objects, and terrain by minimizing the Laplacian deformation energy of an interaction mesh built from corresponding human/robot anatomical keypoints together with sampled object and terrain points. In large-scale scenes, the selection of sampled terrain points is critical. If the sampled terrain points are far from the human, the Laplacian deformation energy may change only marginally even when the local retargeting quality is poor. To address this, we sample not only global terrain points but also additional points in the vicinity of the human. This strategy preserves global geometric proportions while improving local alignment accuracy. We solve for the robot configuration $q_t$ per frame with an SQP-style optimizer, under hard constraints for collision avoidance, joint/velocity limits, and stance-foot anchoring to prevent foot skating.

During retargeting, the human motion may remain collision-free while the retargeted robot still penetrates the terrain due to kinematic mismatch. To improve physical plausibility, we apply a lightweight TSDF-based correction to the robot *global translation*. Specifically, we build a TSDF of the terrain and query the signed distance $d(\mathbf{x})$ at the robot vertices $\mathcal{V}_r^t$, where penetration is indicated by $d(\mathbf{v}) < 0$. Following the feasibility criterion in Eq. 3, we seek an offset $\Delta\mathbf{o}$ such that $\min_{\mathbf{v}\in\mathcal{V}_r^t} d(\mathbf{v} + \Delta\mathbf{o}) \geq \tau_{safety}$, where $\tau_{safety} < 0$ allows a small tolerance. To make the update well-defined, we first compute a unit correction direction $\mathbf{u}$ from the average SDF gradient over penetrated (or near-surface) vertices,

$$\mathbf{u} = \frac{\frac{1}{|\mathcal{M}|}\sum_{\mathbf{v}\in\mathcal{M}}\nabla d(\mathbf{v})}{\left\|\frac{1}{|\mathcal{M}|}\sum_{\mathbf{v}\in\mathcal{M}}\nabla d(\mathbf{v})\right\|_2}, \tag{6}$$

and then parameterize the offset as $\Delta\mathbf{o} = \eta\,\mathbf{u}$. We choose the smallest $\eta \geq 0$ via line search such that the feasibility constraint holds, yielding a collision-free global translation.

## 4. Experiments

### 4.1. Reconstruction Comparison

**Evaluation Setup.**
We evaluate the robustness of our reconstruction pipeline on a subset of the SLOPER4D dataset Dai et al. (2023). Following established evaluation protocols Allshire et al. (2025); Shin et al. (2024); Wang et al. (2024), we report performance on two complementary aspects: (i) *human trajectory reconstruction* and (ii) *scene geometry reconstruction*. For benchmarking, we use a subset of SLOPER4D that includes only sequences where SAM2 tracking—comprising human detection and cross-frame association—is successful. This subset contains two sequences for each activity category: running, walking, and stair ascent/descent.

**Metrics and Baselines.**
For human trajectories, we report World-frame Mean Per Joint Position Error (W-MPJPE) and World-frame Aligned MPJPE (WA-MPJPE). For each sequence, we partition the motion into 100-frame segments. W-MPJPE aligns only the first two frames of each segment to the ground truth, emphasizing global consistency over long horizons, while WA-MPJPE aligns the entire segment and measures local trajectory accuracy over time.

For scene geometry, we report the Chamfer Distance (in meters) between the aligned predicted point cloud and the LiDAR point cloud, restricted to the RGB camera's field of view.

We compare against WHAM Shin et al. (2024), TRAM Wang et al. (2024), and VideoMimic Allshire et al. (2025). WHAM focuses on human motion reconstruction and does not recover the environment (thus Chamfer Distance is not applicable). For fair comparison, all baselines are reproduced using their official implementations.

**Results.**
As summarized in Table 1, our method achieves the best overall performance on both human motion and scene reconstruction. Compared with VideoMimic, we reduce WA-MPJPE from 112.13 to 94.32 ($-15.9\%$) and W-MPJPE from 696.62 to 518.98 ($-25.5\%$), indicating improved local fidelity and substantially better global trajectory stability. For scene geometry, our Chamfer Distance decreases from 0.75 to 0.61 ($-18.7\%$), suggesting more accurate and less noisy terrain reconstruction.

Compared with TRAM, our improvements are more pronounced, especially for geometry: Chamfer Distance drops from 10.66 to 0.61, highlighting the benefit of explicitly modeling and reconstructing the surrounding environment. Overall, these gains validate that our reconstruction pipeline provides higher-quality human–scene inputs, which are critical for downstream contact reasoning and physics-based humanoid learning.
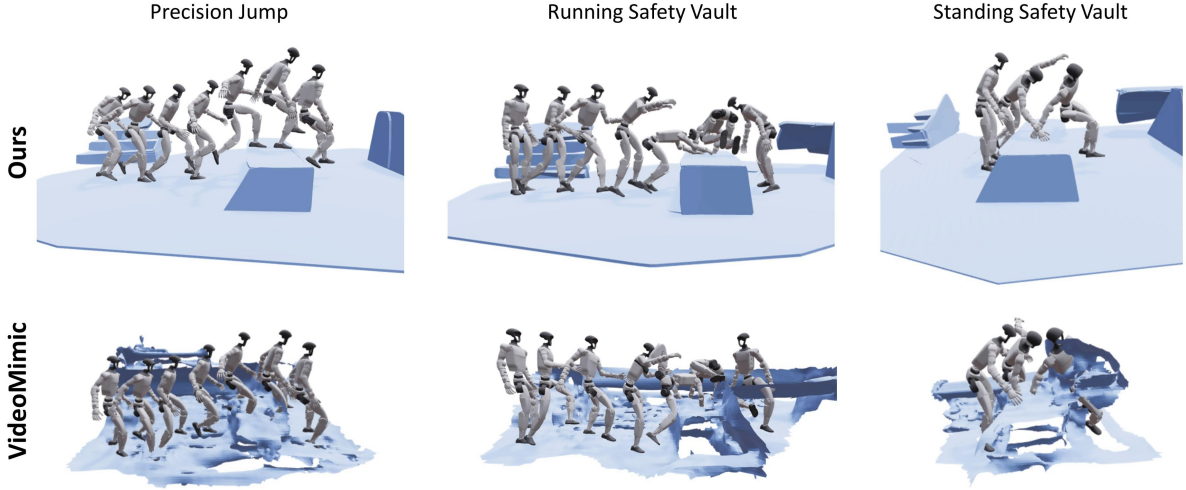
Precision Jump  Running Safety Vault  Standing Safety Vault



Figure 5: **Comparison with VideoMimic.**

| Methods | WA-MPJPE | W-MPJPE | Chamfer Distance |
|---|---|---|---|
| WHAM* Shin et al. (2024) | 189.29 | 1148.49 | – |
| TRAM Wang et al. (2024) | 149.48 | 954.90 | 10.66 |
| VideoMimic Allshire et al. (2025) | 112.13 | 696.62 | 0.75 |
| **Ours** | **94.32** | **518.98** | **0.61** |

Table 1: **Comparison of Reconstruction.** ∗ WHAM does not recover the environment.

## 4.2. Training and Deployment Configuration

Our method leverages high-fidelity human–scene reconstruction to obtain both high-quality kinematic reference motions and watertight scene geometry. Such accurate reconstruction significantly reduces the burden on reward shaping: rather than relying on elaborate humanoid RL reward engineering Li et al. (2025); Weng et al. (2025); Zhao et al. (2025), we adopt a minimal BeyondMimic-style formulation Liao et al. (2025). As summarized in Table 2, our configuration uses only generic tracking terms and standard regularization, yet it is sufficient for robust whole-body tracking with contact-rich scene interactions. Beyond the BeyondMimic-style observation design, we additionally incorporate a global torso position signal during training; during real-robot deployment, we obtain torso position from an optical motion-capture system. We train the policy in IsaacLab Mittal et al. (2025) using asymmetric PPO, where the actor operates on proprioceptive and reference features while the critic has access to privileged scene. Unlike prior interactive humanoid controllers that introduce contact- or task-specific reward heuristics, we find that high-quality reconstruction and reference motions already provide strong implicit supervision.

Training interactive imitation directly from scratch is computationally demanding due to slower simulation caused by complex contact computation. To improve efficiency, we first pre-train a generic whole-body motion tracker Chen et al. (2025); Luo et al. (2025) on ∼50 hours of non-interactive human motion data using the same asymmetric PPO setup. We then fine-tune the policy on scene-interactive motion references. In practice, we adopt lightweight feed-forward architectures for asymmetric PPO: both actor and critic are 4-layer MLPs with hidden dimensions $[3072, 1536, 768, 512]$ and operate with 5-step observation histories and a 5-step future motion horizon. After fine-tuning on scene-interactive data, the resulting policy runs onboard a Unitree G1 robot at 50 Hz using an NVIDIA Jetson Orin. Because the reconstruction pipeline preserves exact real-world terrain geometry, we deploy the robot within the same real-world scenes used for video capture.

| Category | Description |
|---|---|
| **Actor Observation** | Reference Motion: joint position/velocity, future steps torso position/orientation error<br>Proprioception: projected gravity, torso angular velocity, joint position/velocity all with histories<br>Previous Action: executed policy action with histories |
| **Critic Observation** | Reference Motion: joint position/velocity, future steps torso position/orientation error<br>Proprioception: projected gravity, torso linear/angular velocity, body links position/orientation, joint position/velocity with histories<br>Previous Action: executed policy action with histories |
| **Reward** | Anchor Tracking: anchor position/velocity/orientation error<br>Body Tracking: tracking term for body position, orientation, linear and angular velocity<br>Action Rate: penalize rapid action change<br>Soft Joint Limit: penalize limit violation |
| **Termination** | Large anchor position tracking deviation on Z-axis<br>Large anchor orientation tracking deviation<br>Large body position tracking deviation on Z-axis |
| **Domain Randomization** | Randomize robot body material<br>Torso COM: $\pm0.025$m (x), $\pm0.05$m (y), $\pm0.05$m (z)<br>Joint default position: $\pm0.01$ rad<br>Random push: 0.3 m/s, 0.78 rad/s for (1–3) s |
| **Motion Adaptive Sampling** | BeyondMimic-style motion bin adaptive sampling |

Table 2: Training configuration.

## 4.3. Real2Sim2Real Comparison

**Scene Setup.**
We evaluate our real2sim2real pipeline across eight diverse scene-interaction tasks that involve stepping, vaulting, climbing, and parkour-like contact patterns. Specifically, the testing scenes are:

- **walk1**: walk on a flat plane.
- **jump box1 (JB1)**: jump onto a 40 cm box.
- **jump box2 (JB2)**: running single-leg jump onto a 40 cm box and drop down.
- **climb box1 (CB1)**: climb onto a 50 cm box, walk to edge, and descend using single-hand support.
- **climb box2 (CB2)**: side climb onto a 60 cm box.
- **safety vault1 (SV1)**: single-hand safety vault over a 40 cm box.
- **safety vault2 (SV2)**: double-hand safety vault over a 40 cm box.
- **jump climb down1 (JCD1)**: jump onto a 20 cm box, climb onto a 60 cm box, and descend using single-hand support.

(Please refer to the supplementary video for visualizations.)

**Baselines and Metrics.**
We primarily compare against **VideoMimic** Allshire et al. (2025), which similarly provides a monocular real2sim2real pipeline. The comparison is conducted at two levels: (i) *real2sim* via the mean training reward in IsaacLab, and (ii) *sim2real* via real-world success rate (SR). The mean rewards follows the formulation in Table 2 and is measured after 40k PPO iterations with 2048 parallel IsaacLab environments. SR is computed over 10 real-world trials per scene, where a trial is counted as successful if the robot completes the full motion sequence without intervention.

**Terrain and Motion Reconstruction.**
As shown in Fig. 5, VideoMimic struggles to reconstruct terrain geometry due to background clutter and limited camera viewpoints, leading to blurred edges, floating or hollow surfaces, and uneven ground topology. Such
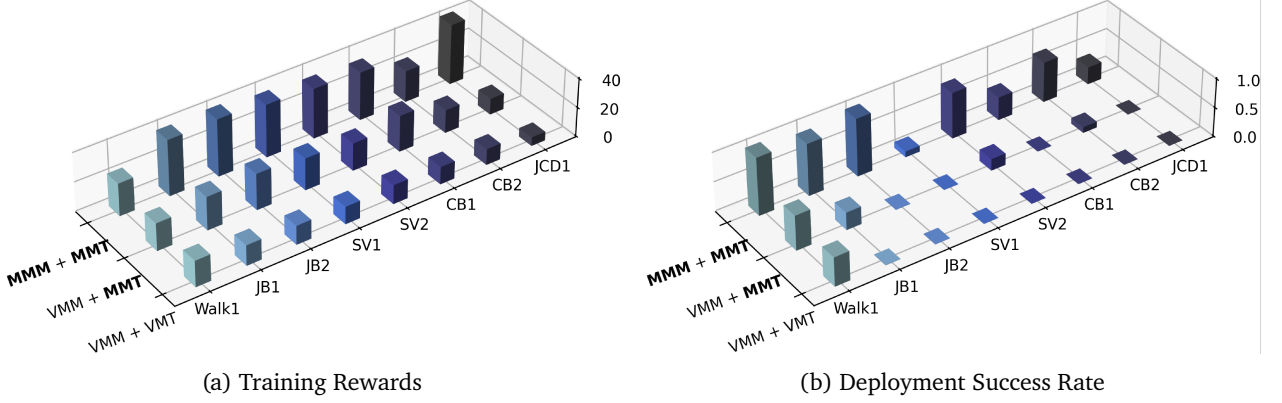
(a) Training Rewards          (b) Deployment Success Rate

Figure 6: Effect of motion and terrain reconstruction on training and deployment performance (**MMM: MeshMimic Motion**; VMM: VideoMimic Motion; **MMT: MeshMimic Terrain**; VMT: VideoMimic Terrain).

artifacts significantly complicate contact reasoning during training and deployment. To decouple motion and terrain effects, we evaluate reconstruction in three configurations: **MMM+MMT** (MeshMimic motion + MeshMimic terrain), **VMM+MMT** (VideoMimic motion + MeshMimic terrain), and **VMM+VMT** (VideoMimic motion + VideoMimic terrain). Quantitative results are shown in Fig. 6.

In terms of *Training reward* (Fig. 6a), VideoMimic motions exhibit frequent foot-in-air artifacts, interpenetrations, and frame-to-frame drift. These errors make it difficult for the policy to track under physical constraints, especially in long-horizon or contact-intensive scenes. As a result, MMM+MMT consistently achieves higher mean rewards than VMM+MMT, with the largest gaps observed in long or multi-contact scenes such as JCD1. Using VMM+VMT further reduces reward due to inaccurate terrain geometry: floating obstacles block humanoid trajectories, hollow geometry causes unexpected falls during contact, and irregular surfaces interfere with adaptive initialization during training. These effects collectively increase motion-tracking difficulty during training.

For *real-world deployment* (Fig. 6b), MMM+MMT achieves stable success rates across most tasks, with failures mainly occurring in highly dynamic settings (e.g., SV1) or long-horizon multi-scene interactions (e.g., CB1, JCD1). In contrast, VMM+MMT fails sim2sim validation (IsaacLab → MuJoCo) on several tasks (JB2, SV1, CB1, JCD1), preventing safe deployment. Even for tasks that pass sim2sim checks (Walk1, JB1, SV2, CB2), deployment remains unstable due to motion drift inherent in VideoMimic reconstructions. With VMT, sim2sim failures become more severe (JB1, JB2, SV1, SV2, CB1, CB2, JCD1), again preventing deployment except in Walk1, where terrain interaction is minimal.

Overall, **MMM+MMT** yields the highest simulation reward and the highest real-world success rate. We attribute this to accurate reconstruction of both motion and terrain, which reduces model mismatch throughout the real2sim2real loop and enables reliable contact reasoning in challenging parkour-like scenes.

**Global Torso Position as an Observation.**
We further investigate whether exposing the policy to global torso position improves sim-to-real deployment. This signal is unavailable from onboard proprioception and must be externally estimated during deployment (e.g., via optical motion capture), thus serving as an exteroceptive cue.

As shown in Table 3, injecting global torso position significantly benefits long-horizon traversal behaviors that require sustained locomotion and multi-contact interactions across obstacles. In particular, we observe success rate improvements of $+20\%$ on JB2, $+20\%$ on CB1, and $+30\%$ on JCD1. These motions span meters and involve transitions between stepping, jumping, and climbing, during which small drift in global pose accumulates into meter-scale deviation. Access to global position mitigates this drift, improving foot placement accuracy and

|         | Walk1 | JB1 | JB2 | SV1 | SV2 | CB1 | CB2 | JCD1 |
|---------|-------|-----|-----|-----|-----|-----|-----|------|
| w/o pos | 1.0   | 1.0 | 0.8 | 0.7 | 1.0 | 0.2 | 1.0 | 0.0  |
| w/ pos  | 1.0   | 0.9 | 1.0 | 0.5 | 0.8 | 0.4 | 0.7 | 0.3  |

Table 3: Effect of adding global torso position to the observation space on sim-to-real deployment success rate.

reducing failure cases where the robot reaches an obstacle "off-phase" or misaligns with the terrain.

Conversely, short-duration yet highly dynamic motions (e.g., SV1, SV2, CB2) show degraded performance, with decreases of $-20\%$, $-20\%$, and $-30\%$ respectively. We observed two contributing factors: (1) fast upper-body and contact transitions introduce intermittent marker occlusions during motion capture, amplifying noise in global position estimates, and (2) these tasks rely more heavily on local agility and contact timing than on accurate global pose. In such regimes, noisy global inputs may destabilize tracking behavior and hinder recovery from mis-steps.

Taken together, these results suggest that global observability is beneficial for long-horizon and path-dependent behaviors, while purely proprioceptive policies remain preferable for short, high-acceleration motions where perception noise dominates the benefit of improved global alignment.

## 5. Conclusion and Future Work

We present a humanoid motion learning framework built purely on RGB video data: without any motion-capture system, markers, or specialized sensing hardware, we directly reconstruct both human motions and the surrounding terrain from videos and transfer them into agile parkour-style skills with robust tracking. This "in-the-wild" data pipeline makes the approach readily applicable to arbitrary outdoor scenarios and significantly lowers the cost of data collection compared to traditional MoCap-based solutions. Through extensive experiments, our method demonstrates stable performance across diverse motions and disturbances, indicating strong potential for deployment on real robotic platforms. As future work, we will push toward vision-based generalized parkour over diverse, previously unseen terrains, and develop a fully closed-loop system that tightly couples perception, planning, and control for long-horizon navigation.

# References

[1] 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 preprint arXiv:2505.03729*, 2025. 5

[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 preprint arXiv:2505.03729*, 2025. 2, 9, 10, 11

[3] Joao Pedro Araujo, Yanjie Ze, Pei Xu, Jiajun Wu, and C Karen Liu. Retargeting matters: General motion retargeting for humanoid motion tracking. *arXiv preprint arXiv:2510.02252*, 2025. 4

[4] Zixuan Chen, Mazeyu Ji, Xuxin Cheng, Xuanbin Peng, Xue Bin Peng, and Xiaolong Wang. Gmt: General motion tracking for humanoid whole-body control. *arXiv preprint arXiv:2506.14770*, 2025. 10

[5] Xuxin Cheng, Yandong Ji, Junming Chen, Ruihan Yang, Ge Yang, and Xiaolong Wang. Expressive whole-body control for humanoid robots. *arXiv preprint arXiv:2402.16796*, 2024. 5

[6] Brian Curless and Marc Levoy. A volumetric method for building complex models from range images. In *Proceedings of the 23rd annual conference on Computer graphics and interactive techniques*, pages 303–312, 1996. 8

[7] 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 *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 682–692, June 2023. 9

[8] Alex Fisher, Ricardo Cannizzaro, Madeleine Cochrane, Chatura Nagahawatte, and Jennifer L Palmer. Colmap: A memory-efficient occupancy grid mapping framework. *Robotics and Autonomous Systems*, 142: 103755, 2021. 4

[9] Jinrui Han, Weiji Xie, Jiakun Zheng, Jiyuan Shi, Weinan Zhang, Ting Xiao, and Chenjia Bai. Kungfubot2: Learning versatile motion skills for humanoid whole-body control. *arXiv:2509.16638*, 2025. 5

[10] 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 preprint arXiv:2406.08858*, 2024. 4, 5

[11] 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 *IEEE/CVF Conf. on Computer Vision and Pattern Recognition (CVPR)*, pages 13274–13285, June 2022. 6

[12] Krishna Murthy Jatavallabhula, Alihusein Kuwajerwala, Qiao Gu, Mohd Omama, Tao Chen, Alaa Maalouf, Shuang Li, Ganesh Iyer, Soroush Saryazdi, Nikhil Keetha, et al. Conceptfusion: Open-set multimodal 3d mapping. *arXiv preprint arXiv:2302.07241*, 2023. 4

[13] Mazeyu Ji, Xuanbin Peng, Fangchen Liu, Jialong Li, Ge Yang, Xuxin Cheng, and Xiaolong Wang. Exbody2: Advanced expressive humanoid whole-body control. *arXiv preprint arXiv:2412.13196*, 2024. 5

[14] Lei Ke, Mingqiao Ye, Martin Danelljan, Yu-Wing Tai, Chi-Keung Tang, Fisher Yu, et al. Segment anything in high quality. *Advances in Neural Information Processing Systems*, 36:29914–29934, 2023. 4

[15] Bernhard Kerbl, Georgios Kopanas, Thomas Leimkühler, and George Drettakis. 3d gaussian splatting for real-time radiance field rendering. *ACM Trans. Graph.*, 42(4):139–1, 2023. 3, 4

[16] Justin Kerr, Chung Min Kim, Ken Goldberg, Angjoo Kanazawa, and Matthew Tancik. Lerf: Language embedded radiance fields. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 19729–19739, 2023. 4

[17] Yanghao Li, Hanzi Mao, Ross B. Girshick, and Kaiming He. Exploring plain vision transformer backbones for object detection. *ArXiv*, abs/2203.16527, 2022. 6

[18] Zhongyu Li, Xue Bin Peng, Pieter Abbeel, Sergey Levine, Glen Berseth, and Koushil Sreenath. Reinforcement learning for versatile, dynamic, and robust bipedal locomotion control. *The International Journal of Robotics Research*, 44(5):840–888, 2025. 10

[19] 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. 5, 10

[20] Zhengyi Luo, Jinkun Cao, Kris Kitani, Weipeng Xu, et al. Perpetual humanoid control for real-time simulated avatars. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 10895–10904, 2023. 4

[21] Zhengyi Luo, Ye Yuan, Tingwu Wang, Chenran Li, Sirui Chen, Fernando Castañeda, Zi-Ang Cao, Jiefeng Li, David Minor, Qingwei Ben, et al. Sonic: Supersizing motion tracking for natural humanoid whole-body control. *arXiv preprint arXiv:2511.07820*, 2025. 5, 10

[22] Naureen Mahmood, Nima Ghorbani, Nikolaus F Troje, Gerard Pons-Moll, and Michael J Black. Amass: Archive of motion capture as surface shapes. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 5442–5451, 2019. 2

[23] Ben Mildenhall, Pratul P Srinivasan, Matthew Tancik, Jonathan T Barron, Ravi Ramamoorthi, and Ren Ng. Nerf: Representing scenes as neural radiance fields for view synthesis. *Communications of the ACM*, 65(1):99–106, 2021. 3, 4

[24] Mayank Mittal, Pascal Roth, James Tigue, Antoine Richard, Octi Zhang, Peter Du, Antonio Serrano-Muñoz, Xinjie Yao, René Zurbrügg, Nikita Rudin, et al. Isaac lab: A gpu-accelerated simulation framework for multi-modal robot learning. *arXiv preprint arXiv:2511.04831*, 2025. 10

[25] Thomas Müller, Alex Evans, Christoph Schied, and Alexander Keller. Instant neural graphics primitives with a multiresolution hash encoding. *ACM transactions on graphics (TOG)*, 41(4):1–15, 2022. 4

[26] Richard A Newcombe, Shahram Izadi, Otmar Hilliges, David Molyneaux, David Kim, Andrew J Davison, Pushmeet Kohi, Jamie Shotton, Steve Hodges, and Andrew Fitzgibbon. Kinectfusion: Real-time dense surface mapping and tracking. In *2011 10th IEEE international symposium on mixed and augmented reality*, pages 127–136. Ieee, 2011. 8

[27] Chaoyi Pan, Changhao Wang, Haozhi Qi, Zixi Liu, Homanga Bharadhwaj, Akash Sharma, Tingfan Wu, Guanya Shi, Jitendra Malik, and Francois Hogan. Spider: Scalable physics-informed dexterous retargeting. *arXiv preprint arXiv:2511.09484*, 2025. 4

[28] 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 *Proceedings IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*, 2019. 6

[29] Xue Bin Peng, Pieter Abbeel, Sergey Levine, and Michiel Van de Panne. Deepmimic: Example-guided deep reinforcement learning of physics-based character skills. *ACM Transactions On Graphics (TOG)*, 37 (4):1–14, 2018. 5

[30] Xue Bin Peng, Yunrong Guo, Lina Halper, Sergey Levine, and Sanja Fidler. Ase: Large-scale reusable adversarial skill embeddings for physically simulated characters. *ACM Transactions On Graphics (TOG)*, 41(4):1–17, 2022. 5

[31] Nikhila Ravi, Valentin Gabeur, Yuan-Ting Hu, Ronghang Hu, Chaitanya Ryali, Tengyu Ma, Haitham Khedr, Roman Rädle, Chloe Rolland, Laura Gustafson, Eric Mintun, Junting Pan, Kalyan Vasudev Alwala, Nicolas Carion, Chao-Yuan Wu, Ross Girshick, Piotr Dollár, and Christoph Feichtenhofer. Sam 2: Segment anything in images and videos. *arXiv preprint arXiv:2408.00714*, 2024. URL https://arxiv.org/abs/2408.00714. 6

[32] Johannes L Schonberger and Jan-Michael Frahm. Structure-from-motion revisited. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 4104–4113, 2016. 4

[33] Soyong Shin, Juyong Kim, Eni Halilaj, and Michael J Black. Wham: Reconstructing world-grounded humans with accurate 3d motion. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 2070–2080, 2024. 9, 10

[34] SAM 3D Team, Xingyu Chen, Fu-Jen Chu, Pierre Gleize, Kevin J Liang, Alexander Sax, Hao Tang, Weiyao Wang, Michelle Guo, Thibaut Hardin, Xiang Li, Aohan Lin, Jiawei Liu, Ziqi Ma, Anushka Sagar, Bowen Song, Xiaodong Wang, Jianing Yang, Bowen Zhang, Piotr Dollár, Georgia Gkioxari, Matt Feiszli, and Jitendra Malik. Sam 3d: 3dfy anything in images. 2025. URL https://arxiv.org/abs/2511.16624. 3, 4, 6

[35] Yifan Wang, Jianjun Zhou, Haoyi Zhu, Wenzheng Chang, Yang Zhou, Zizun Li, Junyi Chen, Jiangmiao Pang, Chunhua Shen, and Tong He. $\pi^3$: Permutation-equivariant visual geometry learning. *arXiv preprint arXiv:2507.13347*, 2025. 6

[36] Yufu Wang, Ziyun Wang, Lingjie Liu, and Kostas Daniilidis. Tram: Global trajectory and motion of 3d humans from in-the-wild videos. In *European Conference on Computer Vision*, pages 467–487. Springer, 2024. 9, 10

[37] Yufu Wang, Yu Sun, Priyanka Patel, Kostas Daniilidis, Michael J Black, and Muhammed Kocabas. Prompthmr: Promptable human mesh recovery. In *Proceedings of the Computer Vision and Pattern Recognition Conference*, pages 1148–1159, 2025. 8

[38] Zihan Wang, Jiashun Wang, Jeff Tan, Yiwen Zhao, Jessica Hodgins, Shubham Tulsiani, and Deva Ramanan. Crisp: Contact-guided real2sim from monocular video with planar scene primitives. *arXiv preprint arXiv:2512.14696*, 2025. 6

[39] 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 preprint arXiv:2509.16757*, 2025. 10

[40] Weiji Xie, Jinrui Han, Jiakun Zheng, Huanyu Li, Xinzhe Liu, Jiyuan Shi, Weinan Zhang, Chenjia Bai, and Xuelong Li. Kungfubot: Physics-based humanoid whole-body control for learning highly-dynamic skills. In *The Thirty-ninth Annual Conference on Neural Information Processing Systems*, 2025. URL https://openreview.net/forum?id=LCPoXt0pzm. 5

[41] Pradyumna Yalandur Muralidhar, Yuxuan Xue, Xianghui Xie, Margaret Kostyrko, and Gerard Pons-Moll. Physic: Physically plausible 3d human-scene interaction and contact from a single image. 2025. 7

[42] Lujie Yang, Xiaoyu Huang, Zhen Wu, Angjoo Kanazawa, Pieter Abbeel, Carmelo Sferrazza, C Karen Liu, Rocky Duan, and Guanya Shi. Omniretarget: Interaction-preserving data generation for humanoid whole-body loco-manipulation and scene interaction. *arXiv preprint arXiv:2509.26633*, 2025. 2, 4, 8

[43] Kangning Yin, Weishuai Zeng, Ke Fan, Minyue Dai, Zirui Wang, Qiang Zhang, Zheng Tian, Jingbo Wang, Jiangmiao Pang, and Weinan Zhang. Unitracker: Learning universal whole-body motion tracker for humanoid robots. *arXiv preprint arXiv:2507.07356*, 2025. 5

[44] Siheng Zhao, Yanjie Ze, Yue Wang, C Karen Liu, Pieter Abbeel, Guanya Shi, and Rocky Duan. Resmimic: From general motion tracking to humanoid whole-body loco-manipulation via residual learning. *arXiv preprint arXiv:2510.05070*, 2025. 10