

# VisualMimic

## Visual Humanoid Loco-Manipulation via Motion Tracking and Generation

Shaofeng Yin\* Yanjie Ze\* Hong-Xing Yu C. Karen Liu<sup>†</sup> Jiajun Wu<sup>†</sup>

\*Equal Contribution <sup>†</sup>Equal Advising

Stanford University



Fig. 1: We present VisualMimic, a visual sim-to-real framework for whole-body humanoid loco-manipulation. Videos are available at [visualmimic.github.io](https://visualmimic.github.io).

**Abstract—** Humanoid loco-manipulation in unstructured environments demands tight integration of egocentric perception and whole-body control. However, existing approaches either depend on external motion capture systems or fail to generalize across diverse tasks. We introduce VisualMimic, a visual sim-to-real framework that unifies egocentric vision with hierarchical whole-body control for humanoid robots. VisualMimic combines a task-agnostic low-level keypoint tracker—trained from human motion data via a teacher-student scheme—with a task-specific high-level policy that generates keypoint commands from visual and proprioceptive input. To ensure stable training, we inject noise into the low-level policy and clip high-level actions using human motion statistics. VisualMimic enables zero-shot transfer of visuomotor policies trained in simulation to real humanoid robots, accomplishing a wide range of loco-manipulation tasks such as box lifting, pushing, football dribbling, and kicking. Beyond controlled laboratory settings, our policies also generalize robustly to outdoor environments.

<sup>1</sup>Work was done during Shaofeng Yin’s internship at Stanford University. Shaofeng is now with Tsinghua University.

Videos are available at: [visualmimic.github.io](https://visualmimic.github.io)

### I. INTRODUCTION

How do humans manage to push a box that is too heavy to move with only their arms? We start with vision perception to localize the box and rely on visual feedback to guide our interaction with the box. To generate sufficient force, we might bend down and push with our hands, lean in with the strength of our arms and shoulders, or even nudge the box forward with our feet. In such cases, every part of the body can be brought into play to accomplish the task. These strategies underscore two fundamental aspects of human loco-manipulation: egocentric visual perception and whole-body dexterity.

Equipping humanoid robots with such human-like object interaction abilities has been a long-standing challenge. Current approaches can be categorized into three main paradigms based on tasks: First, locomotion-focused methods [1], [2]

that excel at terrain traversal but do not address object interaction. Second, approaches that rely on external motion capture systems [3], [4] for object state estimation, restricting their deployment to controlled laboratory environments. Third, vision-based methods for object interaction, which follow two distinct paths: 1) imitation learning approaches [5]–[7] that train visuomotor policies via human demonstrations, which are constrained by the scarcity of large-scale demonstration data and result in limited generalization capabilities; and 2) sim-to-real reinforcement learning (RL) methods [8], [9] that offer greater robustness and generalizability; however, vision-based RL is currently limited to simple environmental interactions such as sitting [8] and stair climbing [8], [9], falling significantly short of human-level object interaction abilities, due to the large exploration and action space of humanoid robots.

We aim to take one step forward on the pathway of sim-to-real RL for visual humanoid-object interaction. To make sim-to-real RL generalize better, we adopt a hierarchical design comprising low-level and high-level policies. In such a hierarchical framework, the task-agnostic low-level policy takes care of balanced control and tracks the command sent by the high-level policy, and the task-specific high-level policy generates simplified tracking commands conditioning on egocentric vision input. This design enables more effective task-specific training. We formulate the command interface as body keypoints (root, hands, feet, head) to ensure both compactness and expressiveness.

To obtain a low-level keypoint tracker that performs human-like behaviors while tracking commands, we curate human motion data and supervise the tracker via motion imitation rewards. However, because keypoint commands alone do not capture the entirety of human motion, we observe that the keypoint tracker can track target keypoints while not perfectly producing human-like behaviors. To address this problem, we adopt a teacher–student training scheme: 1) We first train a motion tracker with full access to current and future whole-body motions, thereby capable of precisely following human reference motions; 2) We then distill this motion tracker into a keypoint tracker that operates on simplified keypoint commands. By doing so, our keypoint tracker captures human motion behaviors while still maintaining a compact command space. Notably, our keypoint tracker is task-agnostic and shared across tasks once trained.

Built upon this general keypoint tracker, we train a high-level keypoint generator via sim-to-real RL. Directly training policies via visual RL significantly slows down the training and leads to non-optimal solutions. Therefore, we also apply a teacher–student scheme: 1) We first train a state-based policy with privileged access to object states, enabling them to solve tasks effectively; 2) We then distill the state-based policy into the visuomotor policy that rely solely on egocentric vision and robot proprioception, making it ready for real-world deployment without external object state estimation. To address the large visual sim-to-real gap (Fig. 8), we apply heavy masking to depth images in simulation, approximating

TABLE I: Comparison of methods of different features.

Method	Whole-Body Dex	Loco-Manipulation	Visual Policy
TWIST [5]	✓	✓	✗
VideoMimic [8]	✗	✗	✓
Hitter [3]	✗	✓	✗
Recipe [10]	✗	✗	✓
HEAD [11]	✗	✗	✓
<b>VisualMimic (Ours)</b>	✓	✓	✓

real-world sensor noise.

Due to the exploration nature of RL, we find that the high-level policy training is not stable when the high-level policies explore the action space that is beyond the human motion space (HMS) present in training motion datasets. We adopt two strategies to alleviate this problem: 1) injecting noise during training the low-level policy to help it adapt to potentially noisy commands from the high-level policy, and 2) clipping actions from the high-level policy to keep them within the feasible HMS.

The resulting framework, VisualMimic, enables us to obtain robust and generalizable visuomotor policies that can zero-shot transfer to the real robot, across a broad range of humanoid loco-manipulation tasks, with relatively simple task-specific reward design and without requiring paired human-object motion data. For real-world experiments (Fig. 4 and Fig. 3), we show that our humanoid robot can 1) lift a 0.5-kilogram box to a height of 1 meter, 2) push a very large box (similar height as the robot and weight 3.8 kilograms) straight and steady with its whole body, 3) dribble a football with the fluency of an experienced player, and 4) kick a box forward with alternating feet. Notably, we also show that our visuomotor policies achieve stable performance in outdoor scenarios, showing strong robustness to real-world variability such as lighting changes and uneven ground.

## II. RELATED WORK

### A. Learning Humanoid Loco-Manipulation

Enabling humanoid robots to perform versatile loco-manipulation in unstructured environments, akin to humans, is a long-standing goal for roboticists. Currently, two main pathways are being explored: (1) Imitation learning on real-world data collected via whole-body teleoperation [5]–[7], [12], [13]. While these methods demonstrate promising task versatility, they remain limited by the scarcity of high-quality data and the difficulty of scaling data collection. (2) Sim-to-real reinforcement learning based on large-scale simulation interaction [3], [8], [9], [14]–[16]. These approaches exhibit strong generalization in specific humanoid motor skills, such as terrain traversal [9], [15], box picking [16], and table tennis [3], but are restricted in task diversity compared to imitation learning. Several works remain confined to simulation; for example, HumanoidBench [14] and SkillBlender [17] adopt hierarchical frameworks similar to ours. However, their policies are often excessively jittery or depend on privileged object states, hindering successful deployment in the real world.

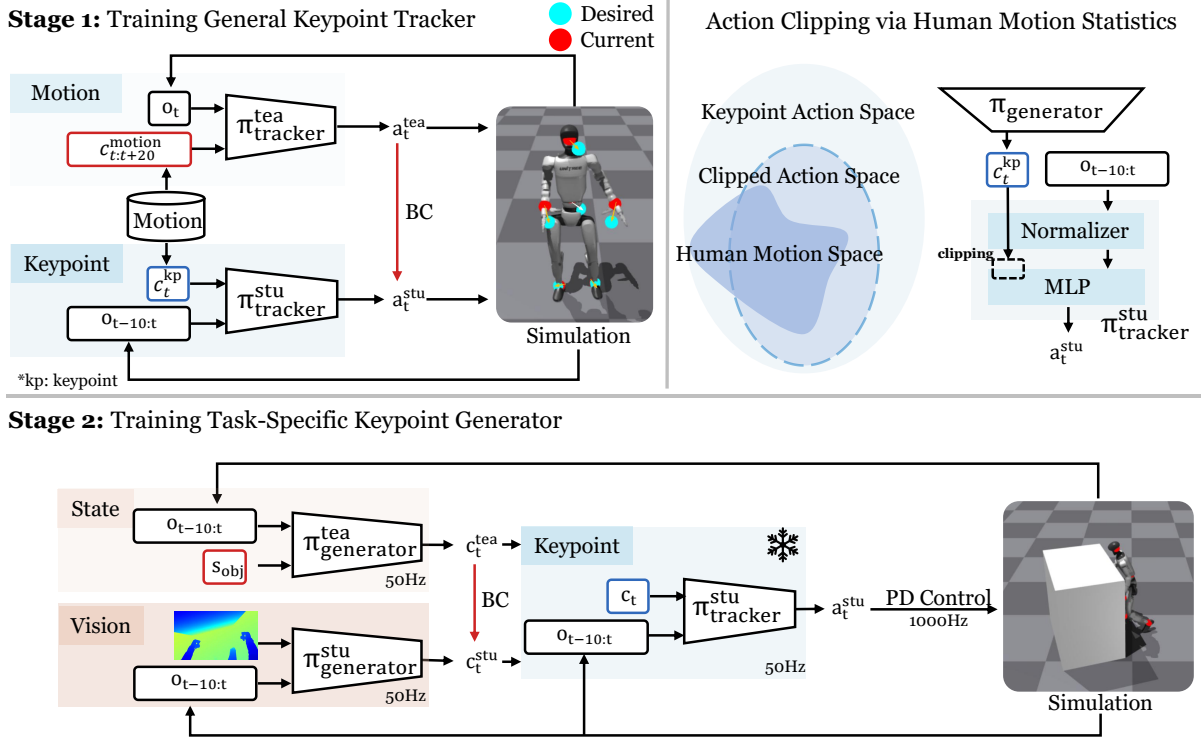


Fig. 2: VisualMimic consists of two training stages: 1) **training a general keypoint tracker**, where a teacher motion tracker is first trained and then distilled into a keypoint tracker with keypoint commands; and (2) **training a task-specific keypoint generator**, where a teacher policy with privileged object states is first trained and then distilled into a visuomotor policy. To ensure stable learning, we compute statistics with human motions and use them to clip high-level actions. Here,  $o_t$  is the proprioceptive observation at time  $t$ ,  $a_t$  is the action, and  $s_{\text{obj}}$  represents the object state.

Previous works on learning visuomotor policies for humanoid robots have typically focused either on upper-body manipulation [10], [18]–[20] or on perceptive locomotion [9], [21], [22]. More recently, VideoMimic [8] introduced a real2sim2real pipeline that enables real robots to perform environment interactions such as sitting, though their interactions remain limited to static settings like the ground or stone chairs. Other efforts, such as PDC [23], have demonstrated promising results in vision-based whole-body loco-manipulation, but only in simulation. In contrast, we propose a sim-to-real framework that enables real-world humanoid robots to achieve versatile object interaction and loco-manipulation through egocentric vision.

### III. METHOD

We propose VisualMimic, a sim-to-real framework for real-world humanoid loco-manipulation by integrating both egocentric visual perception and whole-body dexterity priors within a hierarchical framework (see Fig. 2 for overview). Our approach consists of two main components: 1) a low-level task-agnostic keypoint tracking policy  $\pi_{\text{tracker}}$  that learns whole-body dexterity priors from human motion data, and 2) a high-level task-specific visuomotor policy  $\pi_{\text{generator}}$  that drives the low-level policy based on egocentric visual inputs (Section III-B). Both policies are trained purely via large-scale simulation and zero-shot transfer to the real robot. This hierarchical design facilitates fast adaptation to new tasks,

since only the high-level policy requires training per task. We detail the design of two policies in this section.

#### A. General Keypoint Tracker

Although a keypoint tracking policy can be trained directly, its ability to capture motion is weaker than that of a motion tracking policy due to the much simplified commands, resulting in less human-like behaviors (Fig. 6). We address this problem by designing a two-stage teacher-student distillation pipeline. Specifically, in the teacher training stage, a privileged teacher motion tracker is trained via RL with access to future reference motions. Then, a student keypoint tracker is trained using DAgger [24], relying solely on proprioception and keypoint commands computed from the reference frame at each time step.

*a) Teacher Motion Tracker:* Since the teacher policy is not used during deployment, we provide it with sufficient motion and proprioceptive information so that it can track as accurately as possible. The teacher motion tracker  $\pi_{\text{tea\_tracker}}$  takes as input a sequence of future reference motion frames (spanning 2 seconds) and privileged proprioceptive information (e.g., foot contact forces), which allows it to anticipate upcoming goals and generate smoother motion. We implement  $\pi_{\text{tea\_tracker}}$  as a simple three-layer MLP and optimize it using PPO [25], [26]. Following the reward structure in [5], our reward  $r_{\text{motion}}$  encourages accurate motion tracking while penalizing artifacts such as jitter and foot slippage:





Fig. 3: Our visuomotor policies generalize across diverse space and time, shown on the box-pushing task.



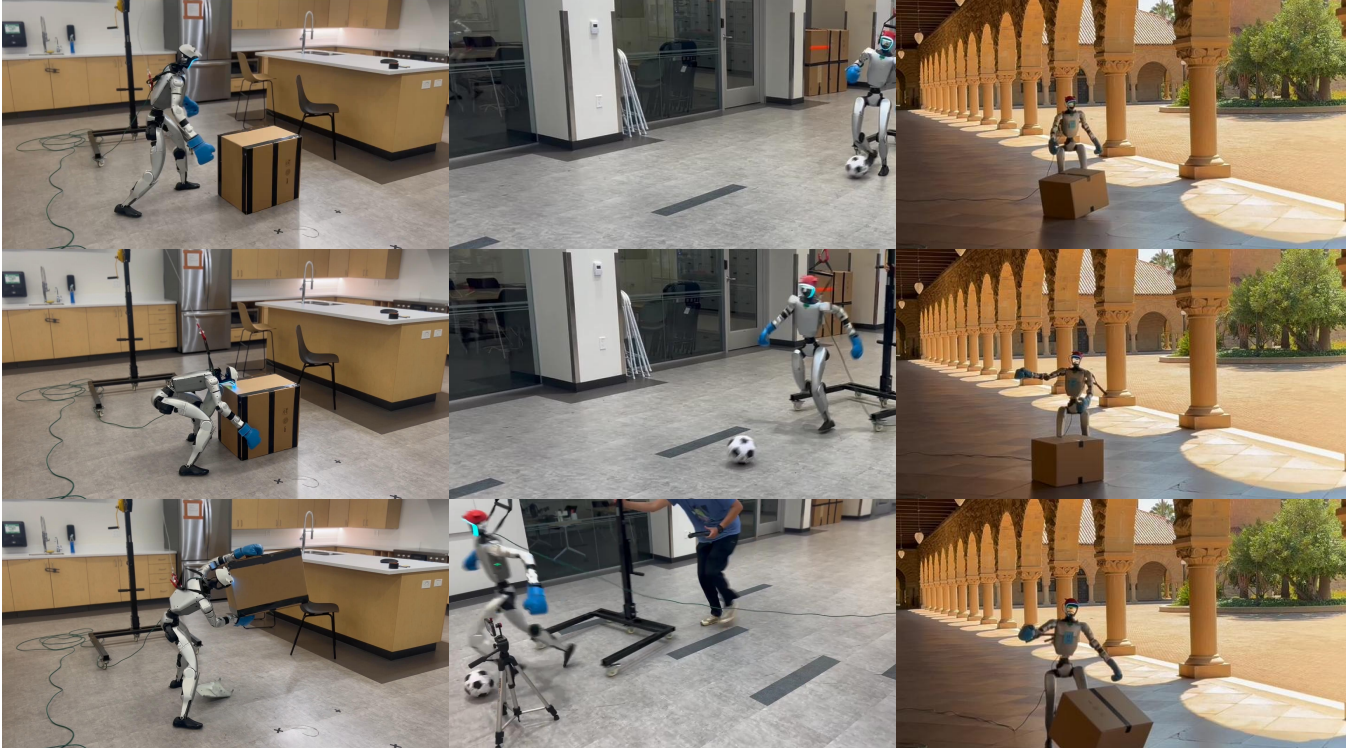


Fig. 4: Real-world deployment of visuomotor policies on a humanoid, showcasing diverse loco-manipulation tasks: Lift Box, Kick Ball, and Kick Box.

$r_{\text{motion}} = r_{\text{track}} + r_{\text{penalty}}$ . Notably, we track robot body positions and root velocities in the world frame. For motion datasets, we use GMR [5], [27] to retarget AMASS [28] and OMOMO [29] into humanoid motions.

*b) Student Keypoint Tracker:* After obtaining  $\pi_{\text{generator}}^{\text{tea}}$ , we distill it into  $\pi_{\text{generator}}^{\text{stu}}$  via Dagger [24], which takes keypoint commands  $c_t^{\text{kp}}$  as input and is real-world deployable. We define

$$c_t^{\text{kp}} = [\Delta p_t, \Delta x_t^1, \dots, \Delta x_t^5], \quad (1)$$

where the root position error is

$$\Delta p_t = p_t^{\text{des}} - p_t^{\text{cur}}, \quad (2)$$

and the keypoint errors (for head, two hands, and two feet) are

$$\Delta x_t^i = (x_t^{i,\text{des}} - p_t^{\text{des}}) - (x_t^{i,\text{cur}} - p_t^{\text{cur}}), \quad i = 1, \dots, 5. \quad (3)$$

Here, “des” and “cur” denote desired (reference) and current trajectories. The keypoint tracker  $\pi_{\text{tracker}}^{\text{stu}}$  relies only on proprioception and the immediate command  $c_t^{\text{kp}}$ . Following the teacher tracker, it is implemented as a three-layer MLP.

### B. Task-Specific Keypoint Generator

Once the low-level keypoint tracker is trained, the next step is to develop a high-level keypoint generator that directs the tracker to perform diverse tasks. Benefiting from the low-level tracker trained with human motion, we only need to focus on guiding the robot toward task completion by

designing a small set of task rewards, without the additional burden of ensuring human-like motion or collecting pairwise human-object interaction data.

Directly training such a keypoint generator with visual RL, however, is highly inefficient (Table III), because crucial information—such as object positions and contact forces—is only partially observable, and incorporating vision into Isaac-Gym further slows down the simulation. To overcome these challenges, we adopt a two-stage approach to train a task-specific keypoint generator. In the first stage, a teacher generator with access to task-relevant object states is trained using PPO [25]. In the second stage, a student keypoint tracker that relies only on depth images for object information is distilled.

*a) Teacher State-based Keypoint Generator:* As discussed in the previous paragraph, the teacher keypoint generator leverages object states to accelerate training. The object state is defined within the environment and concatenated with proprioceptive information to form the input of the state-based keypoint generator. The generator, implemented as a three-layer MLP, is trained using PPO [25], [26] with a task-specific reward function.

We focus on loco-manipulation tasks such as pushing/reaching/kicking objects. Our task rewards are as follows: leftmargin=1.5em

- 1) **Approach** ( $R_{\text{approach}}$ ): encourages contact with target point(s) on target objects. For single-point contact:

$$R_{\text{approach}}(t) = e^{-0.1d(t)}, \quad (4)$$

and for two-point contact such as pushing a box with two hands, we use a harmonic mean for balance:

$$R_{\text{approach}}(t) = \frac{2e^{-0.1d_1(t)}e^{-0.1d_2(t)}}{e^{-0.1d_1(t)} + e^{-0.1d_2(t)}}, \quad (5)$$

where  $d(t)$  (or  $d_1(t), d_2(t)$ ) denotes the distance between a humanoid end-effector (hand or foot) and the target point(s) at time  $t$ .

- 2) **Forward progress** ( $R_{\text{forward}}$ ): rewards new forward motion of the object:

$$R_{\text{forward}}(t) = \tanh\left(10[x_{\text{obj}}(t) - \max_{t' < t} x_{\text{obj}}(t')]\right), \quad (6)$$

where  $x_{\text{obj}}(t)$  denotes the object's position coordinate along the forward direction.

- 3) **Force** ( $R_{\text{force}}$ ): rewarding sufficient force exerted on the object:

$$R_{\text{force}}(t) = e^{-0.1[F_{\text{des}} - F_{\text{obj}}(t)]_+}, \quad (7)$$

where  $F_{\text{obj}}(t)$  is the force exerted on the object and  $F_{\text{des}}$  is the desired force threshold.

Besides, we have the following terms on regularizing the task behavior of the policy: leftmargin=1.5em

- 1) **Look at object** ( $R_{\text{look}}$ ): encourages the robot to face the object:

$$R_{\text{look}}(t) = -(\arccos(\hat{\mathbf{f}}_{\text{body}} \cdot \hat{\mathbf{d}}_{\text{obj}}))^2. \quad (8)$$

$\hat{\mathbf{f}}_{\text{body}}$  denotes the unit vector of the humanoid's facing direction, and  $\hat{\mathbf{d}}_{\text{obj}}$  is the unit vector pointing from the humanoid toward the target object.

- 2) **Drift** ( $R_{\text{drift}}$ ): penalizes lateral deviation:

$$R_{\text{drift}}(t) = -\tanh\left(10[|y_{\text{obj}}(t)| - \max_{t' < t} |y_{\text{obj}}(t')|]\right). \quad (9)$$

$y_{\text{obj}}(t)$  denote the object's position components along the lateral direction.

*b) Student Vision-based Keypoint Generator:* Since object states are unavailable during deployment, the state-based keypoint generator  $\pi_{\text{generator}}^{\text{tea}}$  is distilled into a student keypoint generator that relies solely on visual observations and proprioceptive inputs. Since RGB images have significant sim-to-real gaps, we use depth images as the only visual modality. The depth input is processed by a CNN encoder, whose output is then concatenated with proprioceptive features and passed through an MLP. The student keypoint generator is distilled using DAgger [24].

### C. Clipping Action Space to Human Motion Space

We find that it is hard to maintain the training stability of  $\pi_{\text{generator}}$  even though we use the compact command space, because RL needs heavy exploration during training and the exploration easily goes beyond a feasible space of keypoint commands extracted from human motions. We refer to this feasible space as the **Human Motion Space** (HMS). To alleviate the problem of action exploration beyond the HMS, we propose the following two techniques.

*a) Noised Keypoint Commands for Low-Level Student Training:* To enhance the robustness of the low-level policy and expand the range of Human Motion Space, we inject multiplicative noise into each dimension of the keypoint command during training. Formally, the noised command is defined as  $X_{\text{noised}} = X \cdot \lambda_i$ ,  $\forall i \in \{1, \dots, n\}$ ,  $\lambda_i \sim \mathcal{U}(0.5, 1.5)$ , where  $X$  denotes the original command and  $\lambda_i$  is sampled independently from a uniform distribution. We set the relative noise level to 50%, which is large enough to diversify keypoint commands while preserving motion signals. Empirically, this strategy significantly benefits the subsequent training of the keypoint generator (Fig. 9a).

*b) Action Clip for High-Level Policies:* Besides injecting noise for increasing robustness, we further regularize the output of  $\pi_{\text{generator}}$ . To this end, we first estimate the HMS boundary using the low-level policy input normalizer, and then apply action clip to constrain the high-level policy output within this range. Specifically, each input dimension is modeled as a Gaussian distribution, and the feasible output range for high-level policy is defined as  $\mu \pm 1.64\sigma$ , which covers approximately 90% of the probability mass. The mean  $\mu$  and standard deviation  $\sigma$  are recorded during low-level policy training. Fig. 9b shows that action clip significantly stabilizes the training of the keypoint generator.

### D. Real World Deployment

*a) Vision-Based Sim-to-Real Transfer:* We observe that the depth images from the RealSense camera are highly noisy. To mitigate this, we apply spatial and temporal filters to smooth the real-world depth images. As shown in Fig. 8, even after smoothing, a significant gap remains between simulated and real-world depth images. To address this issue, we apply heavy random masking during training to better approximate real-world visual noise. Specifically, with 20% probability we apply a fixed bottom-left white mask, and with 10% probability each we add up to six independently sampled rectangular masks. These masks are filled with white, black, or gray values, where gray is drawn uniformly from 0 to 1. Each mask covers up to  $30 \times 30$  pixels on an  $80 \times 45$  frame (25% of the image). Without such masking, the robot exhibits unstable behaviors during deployment. We further notice that the angle of RealSense camera on Unitree G1 has slight drift as the neck is not fixed stably. To account for this effect, we apply randomization on the orientation of the robot camera view by up to  $\pm 5^\circ$ .

*b) Safe Real-World Deployment via Binary Commands:* In real-world deployment, it is crucial for the robot to start/pause/end safely during task execution, as simply terminating the program may cause it to fall and get damaged. Therefore, we introduce a binary command signal (0 or 1) that instructs the robot to either pause or execute the task. The robot can freely switch between the two states and always starts in the pause state. We train this behavior with such reward design: when the command is 0, the task reward is disabled, and when the command is 1, the pause reward is disabled. The pause reward corresponds to tracking a stationary standing motion. Both commands are sampled with 50% probability.



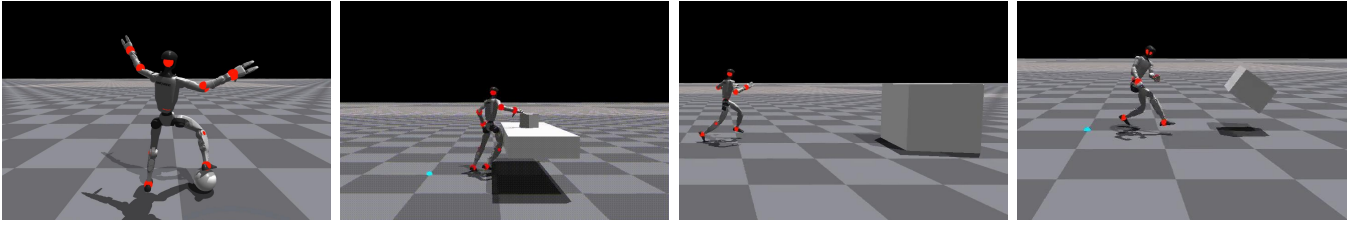


Fig. 5: Visuomotor policies perform diverse loco-manipulation tasks in simulation: from left to right, Balance Ball, Push Cube, Reach Box, Large Kick.

#### IV. EXPERIMENTS AND ANALYSIS

In this section, we perform a series of experiments aimed at addressing the following questions: leftmargin=5pt

- Q1: Does VisualMimic enable effective training of humanoids to perform diverse tasks in a human-like manner?
- Q2: Can policies trained with VisualMimic transfer robustly to the real world?
- Q3: How effectively does our framework demonstrate whole-body dexterity?
- Q4: How well does VisualMimic utilize vision for object interaction?
- Q5: Are design choices of VisualMimic all necessary for the success of the system?

A. *Q1: Does VisualMimic enable effective training of humanoids to perform diverse tasks in a human-like manner?*

We design the following tasks in simulation: leftmargin=5pt

- Push Box:** Push a  $30'' \times 40'' \times 40''$ , 4 kg box (friction 0.5–2.0).
- Kick Box:** Kick a  $15'' \times 20'' \times 20''$ , 0.5 kg box.
- Lift Box:** Lift a  $15'' \times 20'' \times 20''$ , 0.5 kg box.
- Reach Box:** Run to a  $30'' \times 40'' \times 40''$  box.
- Large Kick:** Strongly kick a  $15'' \times 20'' \times 20''$ , 0.5 kg box.
- Kick Ball:** Dribble a football.
- Balance Ball:** Balance with one foot on a football.
- Push Cube:** Push an  $8''$  cube on a tabletop to target.

We evaluate each task using three metrics (TABLE II), with definitions provided in the caption. The choice of metrics follows the reward design of each task. For instance, in Push Box, the reward encourages forward motion along the x-axis while penalizing lateral drift; thus, we report Forward and Drift as metrics. Our framework achieves strong performance and successfully completes all tasks. Our vision-based policy can push a 3.8-kilogram box, comparable in size to the robot, an average of 37 meters per minute. It can also dribble a football an average of 135 meters per minute, with 121 meters forward, indicating both straightness and robustness. This performance stems from our hierarchical design, which decomposes the problem into training a low-level motion tracker and a high-level motion generator.

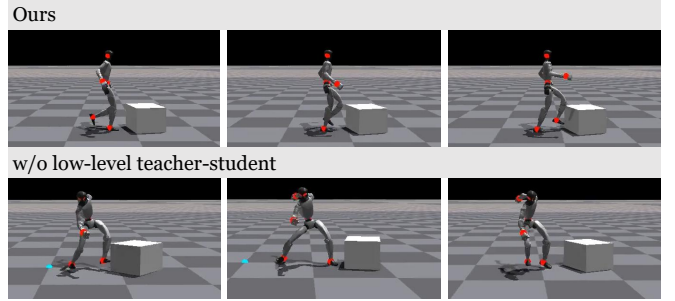


Fig. 6: Box-kicking behaviors. With our teacher–student training (ours, top), the humanoid can mimic human-like motion, while training without it leads to non-human-like motion (bottom).

B. *Q2: Can policies trained with VisualMimic transfer robustly to the real world?*

We deploy our visuomotor policies on the Unitree G1 humanoid robot equipped with its onboard RealSense D435i camera. The policies are evaluated on four tasks (Lift Box, Kick Ball, Kick Box, and Push Box) and all transfer successfully from simulation to the real world (Fig. 4 and Fig. 3). The robot can lift a box from the ground and maintain a stable hold, dribble a football forward with a human-like pose, kick a box in a coordinated manner, and push a large box straight and steadily. Beyond controlled laboratory settings, we also conducted outdoor experiments, in which the policies exhibit robustness to diverse real-world factors (Fig. 3), including fluctuating lighting, irregular ground surfaces, and environmental distractions in the surroundings. These results confirm that our framework enables reliable deployment on physical hardware, allowing the humanoid to perform diverse whole-body loco-manipulation tasks. Supplementary videos are provided.

C. *Q3: How effectively does our framework demonstrate whole-body dexterity?*

Fig. 4 and 5 show that our method enables the robot to interact with objects using its entire body, including the lower limbs. We also observe flexible use of different body parts within the same task (Fig. 7), demonstrating that our method enables adaptive whole-body strategies tailored to environment conditions.

TABLE II: Evaluation of teacher and student policies with and without vision across all tasks in simulation. Results are averaged over 1-minute rollouts across 3 seeds (4096 rollouts for state-based, 512 for vision-based). Metrics include: **Distance** (object Euclidean displacement), **Forward** (object displacement in the intended forward direction), **Drift** (object lateral deviation orthogonal to the forward axis), **Height** (box lift height), **Box Fall Rate** (percentage of drops), **Alive** (time before termination), **Velocity** (forward speed), **Collision Rate** (percentage of collisions), **Force** (average force applied on the ball), **Foot Fall Rate** (percentage of failed foot placements), **Error** (final distance from target), and **Finish Time** (time to task completion; set to episode length if the task is not completed). **The best performing deployable method is highlighted.**

Method	Push Box			Kick Box			Lift Box			Reach Box		
	Distance [m] ↑	Forward [m] ↑	Drift [m] ↓	Distance [m] ↑	Forward [m] ↑	Drift [m] ↓	Height [m] ↑	Box Fall Rate [%] ↓	Alive [s] ↑	Velocity [m/s] ↑	Collision Rate [%] ↓	Alive [s] ↑
teacher	152 ± 36	151 ± 29	13 ± 4	78 ± 3	78 ± 3	0 ± 0	1 ± 0	34 ± 25	38 ± 13	4 ± 0	0 ± 0	60 ± 0
stu w/ vision	<b>37 ± 28</b>	<b>19 ± 15</b>	21 ± 12	<b>55 ± 5</b>	<b>30 ± 3</b>	33 ± 3	<b>1 ± 0</b>	30 ± 23	<b>30 ± 7</b>	<b>4 ± 0</b>	<b>0 ± 0</b>	<b>42 ± 6</b>
stu w/o vision	2 ± 0	2 ± 0	<b>1 ± 0</b>	0 ± 0	0 ± 0	<b>0 ± 0</b>	0 ± 0	<b>15 ± 21</b>	6 ± 4	<b>4 ± 0</b>	<b>0 ± 0</b>	18 ± 6

Method	Large Kick			Kick Ball			Balance Ball			Push Cube (Tabletop)		
	Distance [m] ↑	Forward [m] ↑	Drift [m] ↓	Distance [m] ↑	Forward [m] ↑	Drift [m] ↓	Force [N] ↑	Foot Fall Rate [%] ↓	Alive [s] ↑	Error [cm] ↓	Finish Time [s] ↓	Alive [s] ↑
teacher	8 ± 1	7 ± 1	2 ± 0	189 ± 3	189 ± 3	4 ± 1	21 ± 2	0 ± 0	34 ± 8	9 ± 3	4 ± 1	58 ± 1
stu w/ vision	<b>6 ± 0</b>	<b>6 ± 0</b>	2 ± 0	<b>135 ± 6</b>	<b>121 ± 8</b>	47 ± 12	<b>24 ± 1</b>	<b>0 ± 0</b>	<b>45 ± 7</b>	<b>21 ± 2</b>	<b>20 ± 8</b>	<b>57 ± 0</b>
stu w/o vision	4 ± 0	4 ± 0	<b>1 ± 0</b>	1 ± 0	1 ± 0	<b>0 ± 0</b>	6 ± 0	<b>0 ± 0</b>	5 ± 1	57 ± 22	43 ± 8	51 ± 10

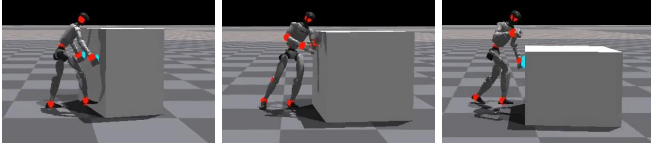


Fig. 7: Whole-body dexterity in Push Box. Policies trained in different conditions exhibit distinct whole-body behaviors: on low-friction ground ( $\mu = 0.5$ ), it bends down to push with both hands; on high friction ( $\mu = 1.5$ ), it leans forward to push with its shoulder for greater force; when the box is lower, it switches to one-handed pushing while the other arm swings with its stride.

To investigate the source of this dexterity, we train a variant of VisualMimic that uses only upper-body keypoints as commands instead of our 6-point design. This variant fails to effectively engage the feet, resulting in poor performance on tasks such as Balance Ball: it struggles to place the foot on top of the ball and instead consistently kicks it away (see supplementary videos).

*D. Q4: How well does VisualMimic utilize vision for object interaction?*

All our tasks rely heavily on robot egocentric vision, where object positions are heavily randomized and robots can only perceive them via vision. To directly assess the role of vision, we train a variant of our policy distilled without visual input (TABLE II). This variant shows a substantial performance drop compared to its vision-enabled counterpart, highlighting the effectiveness of VisualMimic in leveraging vision for robust and adaptive object interaction.

*E. Q5: Are design choices of VisualMimic all necessary for the success of the system?*

To evaluate the contribution of key design components, we conduct extensive ablation studies on three tasks in simulation (TABLE III). We compare against the following variants of VisualMimic: leftmargin=5pt

**w/o noise:** Distill a low-level keypoint tracker without adding noise to keypoint commands, then train the high-level policy on top.

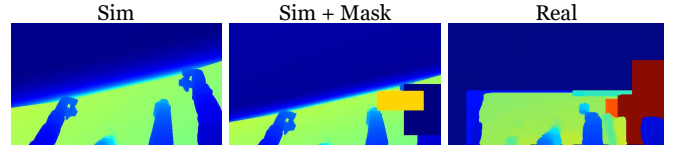


Fig. 8: Egocentric vision of the humanoid robot. In simulation, random masking is applied to approximate real-world noise. Real-world images are processed with spatial and temporal filtering.

**w/o clip:** Use the same low-level keypoint tracker, but remove action clipping in the high-level policy.

**DoF as Interface:** Distill a low-level tracker with a 23-DoF interface, then train the high-level policy on top.

**Local-Frame Tracker:** Train the low-level motion and keypoint tracker with inputs and rewards computed in the local frame, then train the high-level policy on top.

**Visual RL:** Directly train the visuomotor policy with RL for the same number of steps as the combined teacher–student training in the default setup.

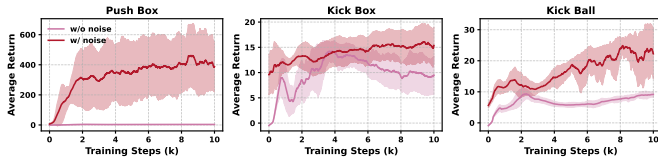
For w/o noise and DoF as Interface, the tracker is distilled from the same teacher policy as in our method. Moreover, for all baselines, only one component is altered while all others are kept identical. Removing multiple components (e.g., DoF interface without noise) leads to even worse performance than single-component variants. All methods share the same reward design and training parameters.

Metrics for the three tasks are reported in TABLE III, with training curves in Fig. 9. It is observed that none of the baselines matches the performance of VisualMimic. We find that removing noise degrades performance in Kick Box and Kick Ball, and causes complete failure in Push Box (Fig. 9a). Eliminating action clipping leads to return collapse in later training stages across all tasks (Fig. 9b), likely because the high-level action space becomes excessively large. Using a DoF interface significantly reduces performance in Push Box and Kick Box, and prevents learning in Kick Ball (Fig. 9c). Training trackers in the local frame also results in substantial drops (TABLE III), as global-frame tracking better reduces drift and facilitates high-level training. Finally, Direct Visual

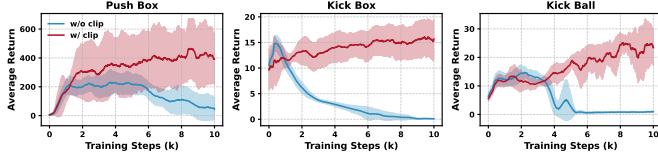


TABLE III: Ablation study of VisualMimic on three tasks in simulation. Experiment settings and metrics follow Table II. **Blind** denotes policies trained without visual input; here, **Ours w/o vision** is the student policy distilled without vision.

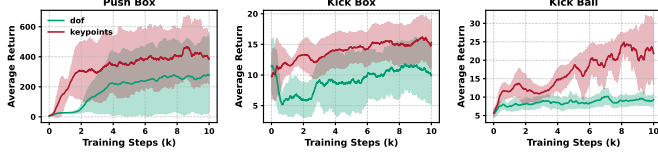
Method	Push Box			Kick Box			Kick Ball		
	Distance [m] ↑	Forward [m] ↑	Drift [m] ↓	Distance [m] ↑	Forward [m] ↑	Drift [m] ↓	Distance [m] ↑	Forward [m] ↑	Drift [m] ↓
<b>State-based</b>									
Ours	<b>152 ± 36</b>	<b>151 ± 29</b>	13 ± 4	<b>78 ± 3</b>	<b>78 ± 3</b>	<b>0 ± 0</b>	<b>189 ± 3</b>	<b>189 ± 3</b>	4 ± 1
w/o noise	2 ± 1	2 ± 1	<b>0 ± 0</b>	30 ± 24	30 ± 20	1 ± 0	136 ± 8	136 ± 7	4 ± 0
w/o clip	68 ± 118	67 ± 94	11 ± 16	3 ± 5	3 ± 4	<b>0 ± 0</b>	12 ± 15	12 ± 12	1 ± 1
DoF as Interface	10 ± 9	8 ± 6	5 ± 3	40 ± 34	40 ± 28	<b>0 ± 0</b>	0 ± 0	0 ± 0	<b>0 ± 0</b>
Local-Frame Tracker	38 ± 27	30 ± 23	16 ± 15	45 ± 7	45 ± 5	1 ± 0	109 ± 23	109 ± 19	7 ± 1
<b>Vision-based</b>									
Ours	<b>37 ± 28</b>	<b>19 ± 15</b>	21 ± 12	<b>55 ± 5</b>	<b>30 ± 3</b>	33 ± 3	<b>135 ± 6</b>	<b>121 ± 8</b>	47 ± 12
w/o noise	2 ± 1	2 ± 1	<b>0 ± 0</b>	25 ± 7	11 ± 4	15 ± 3	86 ± 7	77 ± 7	30 ± 8
w/o clip	10 ± 18	9 ± 12	4 ± 5	6 ± 7	5 ± 3	3 ± 3	1 ± 1	0 ± 1	<b>0 ± 0</b>
DoF as Interface	10 ± 2	6 ± 1	6 ± 1	5 ± 4	1 ± 0	4 ± 3	0 ± 0	0 ± 0	<b>0 ± 0</b>
Local-Frame Tracker	14 ± 11	7 ± 5	8 ± 6	27 ± 15	16 ± 9	15 ± 6	38 ± 17	34 ± 13	12 ± 4
Visual RL	25 ± 16	11 ± 6	16 ± 9	0 ± 0	0 ± 0	<b>0 ± 0</b>	0 ± 0	0 ± 0	<b>0 ± 0</b>
<b>Blind</b>									
Ours w/o vision	2 ± 0	2 ± 0	1 ± 0	0 ± 0	0 ± 0	0 ± 0	1 ± 0	1 ± 0	0 ± 0



(a) Ablation on the noise augmentation for keypoint tracker training.



(b) Ablation on action clip for the high-level policy.



(c) Ablation on the interface design.

Fig. 9: Training curves of our state-based keypoint generator compared with its variants.

RL fails entirely in Kick Box and Kick Ball.

Besides, we find that removing teacher-student training for the low-level tracker, *i.e.*, single-stage RL, produces highly unnatural behaviors during keypoint generator training (Fig. 6), underscoring the importance of the teacher-student framework for human-like keypoint tracking.

## V. CONCLUSIONS AND LIMITATIONS

In this work, we presented VisualMimic, a visual sim-to-real framework that integrates visual perception with whole-body control for humanoid loco-manipulation and object interaction. By combining a high-level keypoint generator with a low-level general keypoint tracker, our approach enables humanoid robots to robustly perform diverse loco-manipulation tasks directly from visual and proprioceptive inputs. Experiments in both simulation and the real world

demonstrate that VisualMimic enables robust deployment from simulation to the real robot.

**Limitations.** While our hierarchical design generalizes across a range of loco-manipulation tasks, more complex interactions involving deformable objects or human collaboration remain unexplored. Besides, though sim-to-real transfer has been effective in our tested scenarios, further scaling to long-horizon tasks and diverse real-world environments may require additional advances in domain randomization and adaptive control. We leave these directions for future research.

## ACKNOWLEDGMENTS

We would like to thank all members of the CogAI group and The Movement Lab from Stanford University for their support. We also thank the Stanford Robotics Center for providing the experiment space. This work is in part supported by Stanford Institute for Human-Centered AI (HAI), Stanford Robotics Center (SRC), ONR MURI N00014-22-1-2740, ONR MURI N00014-24-1-2748, and NSF/FRR 215385.

## REFERENCES

- [1] I. Radosavovic, T. Xiao, B. Zhang, T. Darrell, J. Malik, and K. Sreenath, “Real-world humanoid locomotion with reinforcement learning,” *arXiv:2303.03381*, 2023.
- [2] I. Radosavovic, B. Zhang, B. Shi, J. Rajasegaran, S. Kamat, T. Darrell, K. Sreenath, and J. Malik, “Humanoid locomotion as next token prediction,” *arXiv:2402.19469*, 2024.
- [3] Z. Su, B. Zhang, N. Rahmaman, Y. Gao, Q. Liao, C. Regan, K. Sreenath, and S. S. Sastry, “Hitter: A humanoid table tennis robot via hierarchical planning and learning,” *arXiv preprint arXiv:2508.21043*, 2025.
- [4] H. Xue, C. Pan, Z. Yi, G. Qu, and G. Shi, “Full-order sampling-based mpc for torque-level locomotion control via diffusion-style annealing,” in *2025 IEEE International Conference on Robotics and Automation (ICRA)*. IEEE, 2025, pp. 4974–4981.
- [5] Y. Ze, Z. Chen, J. P. Araújo, Z. ang Cao, X. B. Peng, J. Wu, and C. K. Liu, “Twist: Teleoperated whole-body imitation system,” *arXiv preprint arXiv:2505.02833*, 2025.
- [6] Q. Ben, F. Jia, J. Zeng, J. Dong, D. Lin, and J. Pang, “Homie: Humanoid loco-manipulation with isomorphic exoskeleton cockpit,” *arXiv preprint arXiv:2502.13013*, 2025.
- [7] J. Li, X. Cheng, T. Huang, S. Yang, R. Qiu, and X. Wang, “Amo: Adaptive motion optimization for hyper-dexterous humanoid whole-body control,” *Robotics: Science and Systems 2025*, 2025.

- [8] A. Allshire, H. Choi, J. Zhang, D. McAllister, A. Zhang, C. M. Kim, T. Darrell, P. Abbeel, J. Malik, and A. Kanazawa, "Visual imitation enables contextual humanoid control," *arXiv preprint arXiv:2505.03729*, 2025.
- [9] H. Wang, Z. Wang, J. Ren, Q. Ben, T. Huang, W. Zhang, and J. Pang, "Beamdojo: Learning agile humanoid locomotion on sparse footholds," *arXiv preprint arXiv:2502.10363*, 2025.
- [10] T. Lin, K. Sachdev, L. Fan, J. Malik, and Y. Zhu, "Sim-to-real reinforcement learning for vision-based dexterous manipulation on humanoids," *arXiv preprint arXiv:2502.20396*, 2025.
- [11] S. Chen, Y. Ye, Z.-A. Cao, J. Lew, P. Xu, and C. K. Liu, "Hand-eye autonomous delivery: Learning humanoid navigation, locomotion and reaching," in *9th Annual Conference on Robot Learning*, 2025.
- [12] T. He, Z. Luo, X. He, W. Xiao, C. Zhang, W. Zhang, K. Kitani, C. Liu, and G. Shi, "Omnih2o: Universal and dexterous human-to-humanoid whole-body teleoperation and learning," *arXiv preprint arXiv:2406.08858*, 2024.
- [13] Z. Fu, Q. Zhao, Q. Wu, G. Wetzstein, and C. Finn, "Humanplus: Humanoid shadowing and imitation from humans," in *Conference on Robot Learning (CoRL)*, 2024.
- [14] C. Sferrazza, D.-M. Huang, X. Lin, Y. Lee, and P. Abbeel, "Humanoid-bench: Simulated humanoid benchmark for whole-body locomotion and manipulation," 2024.
- [15] Z. Chen, X. He, Y.-J. Wang, Q. Liao, Y. Ze, Z. Li, S. S. Sastry, J. Wu, K. Sreenath, S. Gupta, and X. B. Peng, "Learning smooth humanoid locomotion through lipschitz-constrained policies," *arXiv preprint arXiv:2410.11825*, 2024.
- [16] J. Dao, H. Duan, and A. Fern, "Sim-to-real learning for humanoid box loco-manipulation," in *2024 IEEE International Conference on Robotics and Automation (ICRA)*. IEEE, 2024, pp. 16 930–16 936.
- [17] Y. Kuang, H. Geng, A. Elhafi, T.-D. Do, P. Abbeel, J. Malik, M. Pavone, and Y. Wang, "Skillblender: Towards versatile humanoid whole-body loco-manipulation via skill blending," *arXiv preprint arXiv:2506.09366*, 2025.
- [18] Y. Ze, Z. Chen, W. Wang, T. Chen, X. He, Y. Yuan, X. B. Peng, and J. Wu, "Generalizable humanoid manipulation with improved 3d diffusion policies," *arXiv preprint arXiv:2410.10803*, 2024.
- [19] X. Cheng, J. Li, S. Yang, G. Yang, and X. Wang, "Open-television: Teleoperation with immersive active visual feedback," *arXiv preprint arXiv:2407.01512*, 2024.
- [20] R.-Z. Qiu, S. Yang, X. Cheng, C. Chawla, J. Li, T. He, G. Yan, D. J. Yoon, R. Hoque, L. Paulsen *et al.*, "Humanoid policy~ human policy," *arXiv preprint arXiv:2503.13441*, 2025.
- [21] J. He, C. Zhang, F. Jenelten, R. Grandia, M. Bächer, and M. Hutter, "Attention-based map encoding for learning generalized legged locomotion," *Science Robotics*, vol. 10, no. 105, 2025.
- [22] J. Long, J. Ren, M. Shi, Z. Wang, T. Huang, P. Luo, and J. Pang, "Learning humanoid locomotion with perceptive internal model," *arXiv preprint arXiv:2411.14386*, 2024.
- [23] Z. Luo, C. Tessler, T. Lin, Y. Yuan, T. He, W. Xiao, Y. Guo, G. Chechik, K. Kitani, L. Fan *et al.*, "Emergent active perception and dexterity of simulated humanoids from visual reinforcement learning," *arXiv preprint arXiv:2505.12278*, 2025.
- [24] S. Ross, G. Gordon, and D. Bagnell, "A reduction of imitation learning and structured prediction to no-regret online learning," in *Proceedings of the fourteenth international conference on artificial intelligence and statistics*. JMLR Workshop and Conference Proceedings, 2011, pp. 627–635.
- [25] J. Schulman, F. Wolski, P. Dhariwal, A. Radford, and O. Klimov, "Proximal policy optimization algorithms," *arXiv preprint arXiv:1707.06347*, 2017.
- [26] N. Rudin, D. Hoeller, P. Reist, and M. Hutter, "Learning to walk in minutes using massively parallel deep reinforcement learning," in *Conference on Robot Learning*. PMLR, 2022, pp. 91–100.
- [27] Y. Ze, J. P. Araújo, J. Wu, and C. K. Liu, "Gmr: General motion retargeting," 2025, gitHub repository. [Online]. Available: <https://github.com/YanjieZe/GMR>
- [28] N. Mahmood, N. Ghorbani, N. F. Troje, G. Pons-Moll, and M. J. Black, "Amass: Archive of motion capture as surface shapes," in *Proceedings of the IEEE/CVF international conference on computer vision*, 2019, pp. 5442–5451.
- [29] J. Li, J. Wu, and C. K. Liu, "Object motion guided human motion synthesis," *ACM Transactions on Graphics (TOG)*, vol. 42, no. 6, pp. 1–11, 2023.