

Physically Plausible Full-Body Hand-Object Interaction Synthesis

Jona Braun¹ Sammy Christen¹ Muhammed Kocabas^{1,2} Emre Aksan^{1†} Otmar Hilliges¹

¹ETH Zurich ²Max Planck Institute for Intelligent Systems, Tübingen

{jona.braun, sammy.christen, muhammed.kocabas, otmar.hilliges}@inf.ethz.ch aksan@google.com

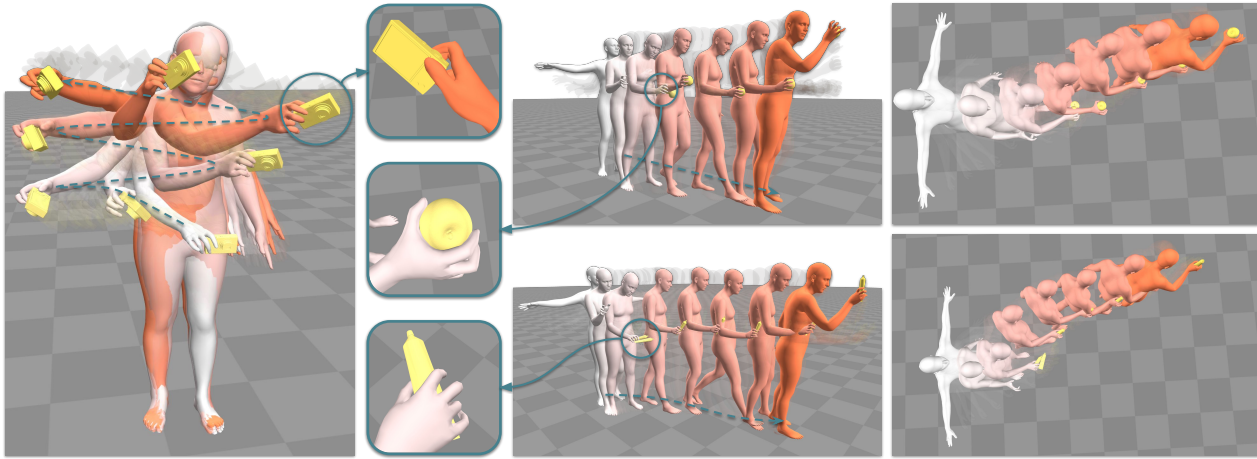


Figure 1. Our method generates physically plausible full-body hand-object interaction sequences. We can synthesize sequences with unseen objects while following a flexibly definable wrist trajectory (left). We can also generate motions of approaching an object, grasping it, then walking to a different location while lifting the object (middle, right). The target trajectories are indicated by the dashed lines.

Abstract

We propose a physics-based method for synthesizing dexterous hand-object interactions in a full-body setting. While recent advancements have addressed specific facets of human-object interactions, a comprehensive physics-based approach remains a challenge. Existing methods often focus on isolated segments of the interaction process and rely on data-driven techniques that may result in artifacts. In contrast, our proposed method embraces reinforcement learning (RL) and physics simulation to mitigate the limitations of data-driven approaches. Through a hierarchical framework, we first learn skill priors for both body and hand movements in a decoupled setting. The generic skill priors learn to decode a latent skill embedding into the motion of the underlying part. A high-level policy then controls hand-object interactions in these pretrained latent spaces, guided by task objectives of grasping and 3D target trajectory following. It is trained using a novel reward function that combines an adversarial style term with a task reward, encouraging natural motions while fulfilling the task incen-

tives. Our method successfully accomplishes the complete interaction task, from approaching an object to grasping and subsequent manipulation. We compare our approach against kinematics-based baselines and show that it leads to more physically plausible motions. Video and code are available at <https://eth-ait.github.io/phys-fullbody-grasp/>.

1. Introduction

Human-object interactions are at the core of our interactions with the physical world. Humans naturally interact with their environment through actions like approaching objects, grasping and manipulating them. The ability to simulate and comprehend these interactions has far-reaching implications in human-computer interaction, robotics, animation and AR/VR.

While recent data-driven works have shown promising results in modeling certain aspects of human-object interactions, a comprehensive, physics-based full-body grasp-

[†]The project was completed after joining Google.

ing approach covering the entire interaction process remains a challenge. Synthesizing dexterous grasps with full-body control is inherently challenging as it requires learning various tasks, namely balancing and moving the body naturally towards the objects, precise finger control, and performing a natural-looking and physically plausible grasp.

Recent works focus on distinct stages of human-object interaction, spanning from the initial approaching phase until grasping [46, 49] to the lifting of objects [9], or even synthesizing the entire sequence [23]. Yet these efforts are primarily data-driven where the intricate physical constraints must be learned from the training data. Such purely data-driven settings can lead to artifacts and unrealistic behaviors due to the inherent limitations of training data such as foot-skating and interpenetration. In contrast, another line of research, physics-based human motion synthesis leverages physics simulation via reinforcement learning (RL) to mitigate limitations of data-driven paradigms. Existing works have either investigated human-object interactions at a larger scale [14, 25] or focused on dexterous hand grasping in an isolated manner [6, 42].

In this paper, we propose the first physics-based method to generate full-body human-object interactions for the entire task of approaching, dexterous grasping and manipulation of objects. By leveraging a physics simulation and reinforcement learning, our method synthesizes natural motions and mitigates physical artifacts, while ensuring that object motions emerge from forces applied by a humanoid agent.

Our method adopts a hierarchical framework, where we first train low-level skill priors and then use these skill priors to learn full-body object interactions. At the core of our approach lies the decoupling of coarse body movement from fine-grained finger control. Specifically, we train separate general-purpose skill priors for the body and hand, decoding latent samples into body and hand movements. This approach ensures that small finger movements are not neglected in a unified training setup. We follow the adversarial training approach to learn these skill priors [34].

To enable full-body object interactions, we build a high-level policy for hand-object interactions that operates in the skill latent spaces. The outputs from this policy are translated into low-level control actions for the physics simulation. The high-level policy can be considered as planning module, leading the entire synthesis process. To guide the training of our high level policy, we propose a novel reward function that combines an adversarial reward to encourage natural motions with a reward to achieve stable grasps. To facilitate the training, we introduce a technique to explicitly condition the policy on 3D target trajectories for the root and wrist positions. This enables the policy to adapt to various scenarios and trajectories during inference.

In this work, we introduce a comprehensive, physics-based approach for the task of full-body grasp synthesis.

Method	Full Body	Physics	Whole Interaction	Dexterous Grasping
ManipNet [54]	×	×	✓	✓
D-Grasp [6]	×	✓	✓	✓
GOAL [46]	✓	×	×	✓
SAGA [49]	✓	×	×	✓
IMOS [9]	✓	×	×	✓
Li <i>et al.</i> [23]	✓	×	✓	✓
Hassan <i>et al.</i> [14]	✓	✓	✓	×
Ours	✓	✓	✓	✓

Table 1. **Method Comparison.** We put our method into context with kinematics-based and physics-based approaches. Our method is the first to achieve physics-based full-body dexterous grasping.

Our method successfully accomplishes the complete interaction task, from approaching (unseen) objects to grasping and subsequent manipulation. We compare our method against the state-of-the-art techniques and present better performance, particularly in physics-based metrics, than the baselines. We further demonstrate the ability to follow diverse and unseen trajectories during inference, showcasing the flexibility and applicability of our method. Our main contributions are as follows:

- A method to generate full-body, dexterous grasping interactions. To the best of our knowledge, this is the first physics-based approach to accomplishing the entire task.
- We propose a two-stage training scheme that decouples dexterous grasping from full-body motion during pre-training and uses joint training during finetuning.
- We compare our method against recent data-driven methods and show that our method produces more physically plausible results.

2. Related Work

We categorize related research into physics-based character control and motion synthesis. Tab. 1 provides an overview of the most related works and ours.

2.1. Physics Based Character Control

Recent research [3, 15, 28, 31–34, 37–39, 51] focuses on using deep reinforcement learning for physics based character control. [28] train a humanoid to catch a tossed ball out of the air and then carry it to a target location. [31] show that incentivizing a policy to follow reference motions through the reward function can generate robust and natural behaviors. In follow-up work, AMP [33] combine adversarial training to imitate reference motions with a task-specific reward. In ASE [34], AMP is scaled to train generalizable skill priors from large motion capture datasets. A high level policy is then trained on the skill prior to fulfill a task objective. In [20], this framework is extended to language conditioned inputs. In contrast to our work, these approaches do not consider finegrained dexterous grasping.

2.2. Motion Reconstruction and Synthesis

Kinematic based The synthesis of human body motion is a well-researched problem in computer vision [1, 2, 10, 16, 17]. Recent work has considered the synthesis of human-scene interaction [4, 13, 18, 22, 43, 48, 56], such as moving a box or sitting on a couch. Contrary to our work, these methods do not consider finegrained hand-object interactions. FLEX [47] jointly optimizes a hand and body pose prior to achieve diverse full-body grasping. Methods that use CVAEs to generate approaching motions for full-body grasps have been proposed [46, 49]. However, the generated motions only model the approaching phase and not the object manipulation phase. On the other hand, a recent work models the object manipulation phase conditioned on language commands [9]. In contrast to these works, we model the full interaction that includes the approaching and manipulation of an object, similar to [23], but employ a physics simulation to increase the physical plausibility of outputs.

Physics-based Recent efforts have been made in leveraging physics simulations for various tasks such as pose estimation [11, 24, 40, 41, 52], human motion synthesis [50], and human-object interaction [5, 6, 14, 25]. Artifacts in pose reconstruction pipelines can for example be corrected by a physics-based policy [11, 24, 25, 52]. [52] use off-the-shelf pose estimation as input to a pretrained imitation learning policy to obtain physically-plausible body motion. [25] extend this by considering indoor scene interactions. [53] learn physically plausible tennis skills from broadcast videos. Most closely related to ours, [14] employ latent skill embeddings from large mocap data [34] and train a high level policy to learn coarse object interaction, such as sitting on a couch or carrying a box. On the other hand, recent works focus on the generation of hand-object interaction sequences in an isolated manner [6, 8, 27, 35, 36]. Approaches often learn dexterous manipulation from full human demonstrations collected via teleoperation [36] or from videos [8, 35]. [27] propose a reward function that incentivizes policies to grasp in the affordance region of objects. [6] propose a reinforcement learning based solution to generate diverse hand-object interactions from sparse reference inputs. However, these approaches either model hand-object interactions but omit the body motion, or focus on the body motion and neglect fine-grained hand-object interactions. In contrast, we generate motions that model full-body hand-object interactions.

3. Task Setting

We model the task of full-body human-object interaction as an RL-problem and leverage a physics simulation for training. We are given an object with global pose $\mathbf{T}_o \in \mathbb{R}^6$ and a human model $\Theta = (\mathbf{t}_b, \theta_b, \theta_h)$, containing the global

translation $\mathbf{t}_b \in \mathbb{R}^3$, the body joint rotations $\theta_b \in \mathbb{R}^{21 \times 6}$ and finger joint rotations $\theta_h \in \mathbb{R}^{16 \times 6}$. We use the continuous 6D representation for rotations [57]. We base the model on the SMPL-X [30] human body model but exclude eye-balls and jaw. Furthermore, we are provided with a hand pose reference Ψ and a target trajectory ξ . The hand-object pose reference captures a single frame of a static hand grasp [6] and is defined as $\Psi = (\bar{\theta}_h, \bar{\mathbf{t}}_h^0, \bar{\mathbf{T}}_o)$, where $\bar{\mathbf{T}}_o$ is the reference object pose, $\bar{\theta}_h$ and $\bar{\mathbf{t}}_h^0$ indicate the target wrist joint rotations and translation, respectively. The target trajectory contains n global target body and wrist 3D positions $\xi = [\mathbf{t}_b^i, \mathbf{t}_h^i]_{i=1}^n$. The goal of the task is to generate an output sequence of human and object poses $[\Theta^t, \mathbf{T}_o^t]_{t=1}^T$ over horizon T . We split the task in two phases; in the first phase, the human character has to walk to the surface with the object and reach a grasp on the object. In the second phase, it has to manipulate the object by consecutively reaching the targets in the trajectory ξ .

3.1. Simulation Environment

In the following we describe the environment of the physics simulation in which we train our human character. We generate a controllable human body model following [52]. It contains 57 DoF actuators for the body joints and 48 DoF actuators for the fingers, totaling 105 DoF. The root of the human (i.e., global 6DoF translation and orientation) is not actuated and changes according to the control of the other body joints. To reduce the computational complexity we approximate the collision geometries of the rigid body meshes with the exception of the ankles and feet. We focus on right-hand grasping and thus omit the left hand’s fingers. We decimate all the object meshes to increase simulation speed. We use proportional derivative (PD) controllers to compute the torques τ to actuate the joints:

$$\begin{aligned} \tau &= k_p \circ (\hat{\theta} - \theta) - k_d \dot{\theta} \\ \hat{\theta} &= \theta_{\text{ref}} + k_s \mathbf{a} \end{aligned} \tag{1}$$

where $\hat{\theta}$ indicate the target joint rotations, θ the current joint rotations, $\dot{\theta}$ the velocity and k_p, k_d, k_s the gains. The target comprises the reference pose θ_{ref} and residual actions \mathbf{a} , which are predicted by our policies. The reference pose θ_{ref} equals the current pose θ_h for the finger control and the center between the joint limits for the body joint control. The state space of the simulation is given by $\mathbf{s} = (\Theta, \dot{\Theta}, \mathbf{T}_o, \dot{\mathbf{T}}_o, \mathbf{f})$, which contains the human pose and velocity information, the object pose and velocity, and the net contact force $\mathbf{f} \in \mathbb{R}^{39 \times 1}$ acting on the human body joints, the object, and the table surface. See supp. material for more details about the simulation environment.

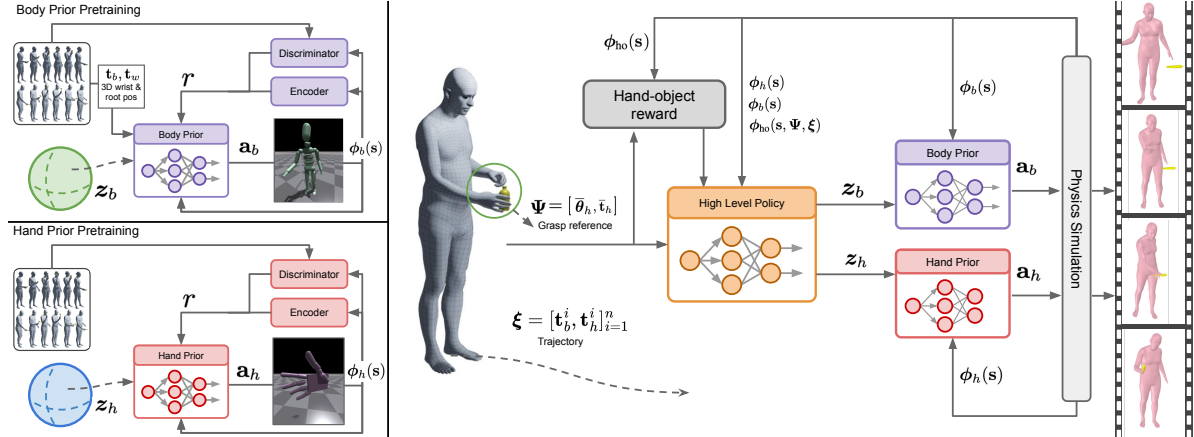


Figure 2. **Method Overview.** Given a hand-object pose reference Ψ and a root and wrist target trajectory ξ , our method synthesizes a human approaching, grasping an unseen object and following a trajectory with it. The training procedure involves learning skill priors, which we decouple into a body prior and a hand prior. We train the priors from a motion capture dataset using adversarial training. The right-hand prior policy $\pi_h(\mathbf{a}_h | \phi_h(s), \mathbf{z}_h)$ predicts the right-hand target joint angles \mathbf{a}_h . The body prior policy $\pi_b(\mathbf{a}_b | \phi_b(s), \mathbf{z}_b, \mathbf{t}_b, \mathbf{t}_w)$ predicts target joint angles of the body \mathbf{a}_b . PD-controllers compute the necessary torques to drive the joints to the predicted target angles in the physics simulation. Instead of directly conditioning the policies on the physics simulation state s , we define feature extraction functions $\phi(\cdot)$ for each policy. A latent vector \mathbf{z} is used to represent the skill space. The body prior is additionally conditioned on a root and wrist 3D target positions \mathbf{t}_b and \mathbf{t}_h . To train a hand-object interaction policy $\pi_{ho}(\mathbf{z}_b, \mathbf{z}_h, \mathbf{t}_b, \mathbf{t}_h | \phi_h(s), \phi_b(s), \phi_{ho}(s, \Psi, \xi))$, we predict the latent vectors \mathbf{z}_b and \mathbf{z}_h of the body and hand prior, respectively. We predict position targets $\bar{\mathbf{t}}_b$ and $\bar{\mathbf{t}}_h$ for the root and wrist as an auxiliary objective. We define a task reward function based on the physics simulation state s , the static hand pose reference Ψ , and the trajectory ξ .

3.2. Reinforcement Learning

We follow [44] and model RL as a Markov Decision Process (MDP) defined by a 6-Tuple $\mathcal{M} = (\mathcal{S}, \mathcal{A}, \mathcal{T}, \mathcal{R}, \mu, \gamma)$, where \mathcal{S} is the state and \mathcal{A} the action space. The deterministic transition function $\mathcal{T} : \mathcal{S} \times \mathcal{A} \mapsto \mathcal{S}$ maps from a state-action pair to the next state and the reward function $\mathcal{R} : \mathcal{S} \times \mathcal{A} \mapsto \mathbb{R}$ maps to a scalar value. The first state s_0 is determined by the initial state distribution $\mu : \Delta(\mathcal{S}) \mapsto \mathcal{S}$. Finally, $\gamma \in [0, 1]$ defines the discount factor of future rewards. We define a parametric policy $\pi : \mathcal{S} \mapsto \Delta(\mathcal{A})$ that maps to a distribution over actions given a state. We aim to optimize the policy such that it maximises the expected discounted reward $\mathbb{E}_{\mathbf{a}_t \sim \pi_\theta(\cdot | s_t), s_0 \sim \mu} \left[\sum_{t=0}^{T-1} \gamma^t \mathcal{R}(s_t, \mathbf{a}_t) \right]$ where $s_{t+1} = \mathcal{T}(s_t, \mathbf{a}_t)$ and T is the horizon.

4. Full-Body Grasp Motion Synthesis

Our framework is inspired by ASE [34] and depicted in Fig. 2. Therefore, we leverage a hierarchical framework. First, we train low-level priors that represents diverse motion skills from motion capture data. Thereafter, we train a high-level policy, dubbed hand-object interaction policy, that predicts actions in the latent spaces of the priors to achieve a high-level objective. In our setting, the objective is to approach the object, grasp it and move it according to a specified wrist and root trajectory. We now first explain how we train physics-based body and hand priors and then describe our hand-object interaction policy training.

4.1. Pre-Training of Body and Hand Priors

In our approach, we decouple the training of the body prior and the hand prior. Crucially, this prevents mode collapse and allows learning coarse body movements and fine-grained finger control. Each prior is represented by a policy $\pi(\mathbf{a} | \phi(s), \mathbf{z})$, which is conditioned on features extracted from the physics simulation’s state $\phi(s)$ and a latent skill vector $\mathbf{z} \sim p(\mathbf{z})$. We combine a motion imitation objective and an unsupervised skill discovery objective [34] to train these priors. The motion imitation objective incentivises the policy to perform motions that are similar as depicted in the reference motion. It is optimized by training a discriminator to differentiate between motions sampled from the reference motion capture data and motions generated by the humanoid character. The skill discovery objective promotes the policy to learn a meaningful latent skill space which allows a high-level policy to reuse the learned skills. Thus, the reward function is defined as follows:

$$r = -\log(1 - D(\phi(s), \phi(s'))) + \beta \log q(\mathbf{z}_t | \phi(s), \phi(s')), \quad (2)$$

where D indicates the discriminator and q is an encoder trained with the objective to recover the latent skill vector \mathbf{z} from a tuple of features $(\phi(s), \phi(s'))$ from the simulation state s and the consecutive state s' .

Hand Prior The hand prior is a policy $\pi_h(\mathbf{a}_h \mid \phi_h(\mathbf{s}), \mathbf{z}_h)$ that controls the wrist and the finger joints via the actions \mathbf{a}_h . It is conditioned on the latent skill vector \mathbf{z}_h and the right-hand features $\phi_h(\mathbf{s}) = (\boldsymbol{\theta}_h, \dot{\boldsymbol{\theta}}_h, \mathbf{x}_h)$, where $\boldsymbol{\theta}_h \in \mathbb{R}^{16 \times 6}$ and $\dot{\boldsymbol{\theta}}_h \in \mathbb{R}^{16 \times 6}$ indicate the local hand joint rotations (except for the global wrist joint orientation) and their angular velocities, and $\mathbf{x}_h \in \mathbb{R}^{16 \times 3}$ are the wrist-relative 3D finger joint positions. To train the hand prior, we detach the hand from the body and fix its global position in space. For training, we use the reward function in Eq. (2) with hand-state tuples $(\phi_h(\mathbf{s}), \phi_h(\mathbf{s}'))$.

Body Prior We extend our body prior setting to a goal-conditioned approach by explicitly considering the 3D target positions of the root \mathbf{t}_b and the wrist \mathbf{t}_h as conditional variables. Our body-prior policy $\pi_b(\mathbf{a}_b \mid \phi_b(\mathbf{s}), \mathbf{z}_b, \mathbf{t}_b, \mathbf{t}_h)$ controls all body joints except the hands. Similarly, the body encoder q_b is conditioned on the target positions such that $q_b(\mathbf{z}_b \mid \phi_b(\mathbf{s}), \phi_b(\mathbf{s}'), \mathbf{t}_b, \mathbf{t}_w)$. We further leverage this additional information by introducing an auxiliary reward on the target positions during high-level policy training (see Section 4.2). The benefits of including the root and wrist targets in the conditional variables are twofold. First, both the policy and the encoder gain spatial awareness, reducing ambiguity and yielding better planning. Second, this formulation allows us to control generated motion at inference time, e.g., walking to a target root position or moving the right wrist to a target position.

The body-state features are defined as $\phi_b(\mathbf{s}) = (\boldsymbol{\theta}_b, \dot{\boldsymbol{\theta}}_b, \mathbf{x}_b, \dot{\mathbf{x}}_b, \mathbf{h}_b, \dot{\mathbf{t}}_b)$. The terms $\boldsymbol{\theta}_b$ and $\dot{\boldsymbol{\theta}}_b$ indicate the root-relative body joint rotations and their velocities (except for the global root joint orientation and velocity). \mathbf{x}_b and $\dot{\mathbf{x}}_b$ are 3D joint positions and their velocities (excluding the root). \mathbf{h}_b is the root’s height (e.g., the value in z-direction according to our preprocessing) and $\dot{\mathbf{t}}_b$ is the root’s linear velocity. All the features except the root height and root orientation are in the root-relative coordinate-frame. The body-state features for the discriminator are a subset of the policy features $\phi_b(\mathbf{s})$, similar to [34]. For training, we use the reward function in Eq. (2) with body-state tuples, $(\phi_b(\mathbf{s}), \phi_b(\mathbf{s}'))$. See supp. material for more details.

4.2. Training of Hand-Object Interaction Policy

We leverage the body and hand prior to train a hand-object interaction policy. The policy $\pi_{ho}(\mathbf{z}_b, \mathbf{z}_h, \dot{\mathbf{t}}_b, \dot{\mathbf{t}}_h \mid \phi_b(\mathbf{s}), \phi_h(\mathbf{s}), \phi_{ho}(\mathbf{s}), \Psi, \xi)$ is conditioned on both hand and body features $\phi_b(\mathbf{s}), \phi_h(\mathbf{s})$ and task-relevant features $\phi_{ho}(\mathbf{s}, \Psi, \xi)$ (see below). It predicts the latent vectors \mathbf{z}_b and \mathbf{z}_h of both the body and hand prior. These latent vectors are then passed to the policies which yield output actions that are applied to the human body model. Additionally, we predict position targets $\dot{\mathbf{t}}_b$ and $\dot{\mathbf{t}}_h$ for the root and wrist as a training scheme which we dub target guidance.

Hand-Object Features The features $\phi_{ho}(\mathbf{s})$ represent the task-relevant information that is required for grasping the object and following a target trajectory:

$$\phi_{ho}(\mathbf{s}, \Psi, \xi) = (\mathbf{T}_o, \dot{\mathbf{T}}_o, \mathbf{g}_x, \mathbf{g}_\theta, \mathbf{g}_c, \mathbf{g}_\xi, \mathbf{f}_h, \mathbf{x}_{\text{tab}}, \eta). \quad (3)$$

The 6D root-relative object pose and its velocity are given by \mathbf{T}_o and $\dot{\mathbf{T}}_o$. The terms $\mathbf{g}_x, \mathbf{g}_\theta$, and \mathbf{g}_c are features computed from the static hand pose reference Ψ (see Section 3) to measure the distance between the current hand pose and the target hand pose:

$$\mathbf{g}_x = \bar{\mathbf{x}}_h - \mathbf{x}_h \quad \mathbf{g}_\theta = \bar{\boldsymbol{\theta}}_h \ominus \boldsymbol{\theta}_h \quad \mathbf{g}_c = (\bar{\mathbf{c}}_h, \bar{\mathbf{c}}_h \ominus \mathbf{c}_h) \quad (4)$$

The distance between the 3D joint positions of the reference pose $\bar{\mathbf{x}}_h$ and the current pose \mathbf{x}_h in root-relative frame is given by \mathbf{g}_x . The 6D rotational difference between the reference hand pose $\bar{\boldsymbol{\theta}}_h$ and the current hand pose $\boldsymbol{\theta}_h$ is defined by \mathbf{g}_θ . Similarly, \mathbf{g}_c is a tuple containing the contact targets and the distance between the target and the current contacts. It is a vector with binary values indicating whether a target contact is achieved or not. The target 3D joint positions $\bar{\mathbf{x}}_h$ and the target contacts $\bar{\mathbf{c}}_h$ are computed from the hand pose reference Ψ . Note that contacts in our context are on a per-joint basis.

Similarly, to guide the human character along a given trajectory, it is provided with the distance to the next waypoints on the trajectory \mathbf{g}_ξ :

$$\mathbf{g}_\xi = (\mathbf{t}_b^i - \mathbf{t}_b, \mathbf{t}_h^i - \mathbf{t}_h), \quad (5)$$

where \mathbf{t}_b^i and \mathbf{t}_h^i are the next root and wrist targets to achieve. Once a target has been reached, the next one is sampled from the trajectory ξ .

Lastly, \mathbf{f}_h is the vector describing the net forces acting on the hand joints, the object, and the table surface (see Section 3.1). The term \mathbf{x}_{tab} is the distance between the 3D wrist joint and the table. The phase variable $\eta \in [0, 1]$ depicts the progress of the task. We provide more details on the hand-object state features $\phi_{ho}(\mathbf{s})$, in supp. material.

Hand-Object Reward Function To guide the policy to grasp the object and follow the trajectory ξ , we define the following hand-object reward function:

$$r_{HO} = w_T r_T(\mathbf{s}, \mathbf{a}) + w_S r_S(\mathbf{s}), \quad (6)$$

where r_T and r_S indicate the task and style reward with weights w_T and w_S , respectively.

The task reward r_T incentivizes the policy to achieve a stable grasp on the object and follow the target trajectory:

$$r_T = r_x + r_\theta + r_c + r_\xi + r_{\text{reg}}, \quad (7)$$

where the terms r_x, r_θ, r_c , and r_ξ are position, orientation, contact and trajectory rewards, respectively. These rewards

are computed by taking the norm of the distance features introduced in Eq. (4) and Eq. (5). Lastly, r_{reg} indicates a regularization reward on the predicted actions. Details on the reward function are provided in the supp. material.

We introduce a style reward r_S to achieve more plausible and natural motions. It extends the discriminator-based style reward of [34] for the hand. Specifically, we use the discriminator predictions for the hand and body such that

$$r_S = -\log(1 - D_b(\phi_b(\mathbf{s}), \phi_b(\mathbf{s}')) - \log(1 - D_h(\phi_h(\mathbf{s}), \phi_h(\mathbf{s}'))). \quad (8)$$

Target Guidance We introduce target guidance to allow the policy to be robust and flexibly follow the given targets $\xi = [\mathbf{t}_b^i, \mathbf{t}_h^i]_{i=1}^n$ for the root and wrist joints. During training, we alternate the target trajectory ξ between the ground-truth $(\bar{\mathbf{t}}_b, \bar{\mathbf{t}}_h)$ and the predicted targets $(\tilde{\mathbf{t}}_b, \tilde{\mathbf{t}}_h)$ and regularize the training with an auxiliary objective:

$$\mathcal{L}_\xi = \|\bar{\xi} - \tilde{\xi}\|_2^2, \quad (9)$$

where $\tilde{\xi}$ is the target prediction and $\bar{\xi}$ is the ground truth trajectory. The loss measures the Euclidean distance between the two terms. Note that target guidance is applied only after the object has been grasped.

4.3. Implementation Details

We follow the actor-critic framework [44] and implement our skill priors with 4-layer MLP networks using [1024, 1024, 512, 512] units and ReLU activations after every layer. In the actor network, we use a Gaussian output model with a constant variance and predict only the mean. The discriminators and encoders share the first 3 linear layers [1024, 1024, 512] with separate final layers. The high-level hand-object policy π_{ho} is implemented with a 3-layer MLP and a Gaussian output model with constant variance where the final layer predicts the mean. For training, we use the Adam optimizer [21] with a learning rate of $2e-5$ and a discount factor γ of 0.99. We implement our method in PyTorch [29]. We use Isaac Gym [26] as physics simulation. It runs at 120Hz while the policies are sampled at 30Hz. Further details can be found in the supp. material.

5. Experiments

We first describe the data and experimental details in Sections 5.1 and 5.2. Section 5.3 presents our main evaluations, consisting of quantitative and qualitative comparisons against the baselines. Lastly, in Section 5.4, we provide an ablation to highlight the contributions of our method.

5.1. Data

We train and evaluate our model using the GRAB dataset [45] where we follow the right-handed grasp setting as in

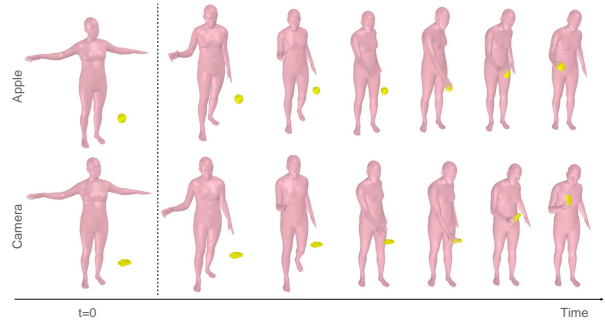


Figure 3. **Qualitative Results.** Each row shows a motion sequence generated by our model for an unseen object from the test set.

the prior works [23, 46, 49]. We combine the object test-split from GOAL [46] and the subject test-split from IMoS [9]. Hence, our training set contains all sequences from subjects S1-S9 and the object-split of GOAL. We then evaluate on both the GOAL and IMoS test sets.

Our humanoid character in the physics simulation is based on the neutral SMPL-X model. Hence, we convert the subject-specific GRAB reference motions to the neutral model. This preprocessing involves aligning the feet with the ground and the object with the hand. We provide more details on the preprocessing in the supp. material.

5.2. Experimental Details

During training of the hand-object interaction policy, we initialize the character at a random frame of the approaching phase sampled from a GRAB reference clip. The object and table are initialized according to the hand pose reference. We use a two-stage training procedure. First, we fix the object to its surface, such that the character can learn to approach and initiate a stable grasp on the object without the risk of moving or dropping the object. In the second stage, the object is non-stationary such that the policy learns to lift and follow the trajectory. To avoid overfitting, we add random noise to the hand pose reference, the target trajectory, the initial object position and rotation around the yaw axis. The noise applied to the object position is also added to the table position to prevent interpenetration of the object.

The one-to-one correspondence between the neutral SMPL-X model and our humanoid in the physics simulation enables a direct conversion between the two. Hence, we are able to run evaluations in the SMPL-X parameter space (except for the grasping success and the TTR metric, see Section 5.2.2) and compare our method against the kinematics-based approaches. At evaluation time, the humanoid agent is always initialized in T-pose and its root is set to the root of the initial test frame. Finally, we apply Gaussian smoothing to the output motion as a post-processing step. We find that the smoothing operation marginally improves the performance. Our model’s performance without the smoothing operation is reported in supp. material.

Method	Success (\uparrow)	GD [mm] (\downarrow)	FS [%](\downarrow)	IV [cm^3] (\downarrow)	ID [mm] (\downarrow)	TTR \uparrow	CR
Approaching							
Ground-truth	0.29	5.9	7.7	1.66	4.4	-	0.111
GOAL [46]	0.13	6.8	14.4	1.97	5.3	-	0.128
Ours	0.79	2.1/0.0*	5.7	0.11	1.1	-	0.026
Manipulation							
Ground-truth (S10 test set)	0.22	6.1	2.9	2.75	4.9	-	0.112
IMOS [9] (S10 test set)	0.20	16.0	8.0	5.07	6.8	-	0.057
Ours (S10 test set)	0.64	1.8/0.0*	0.9	0.22	2.7	0.65	0.053
Ours (GOAL test set)	0.79	1.9/0.0*	1.2	0.18	2.9	0.85	0.055

Table 2. **Evaluation.** We compare our method against the relevant baselines on approaching until grasping and manipulation after grasping. In both settings, we find that our method achieves better performance across all of the metrics. We also provide the metrics for the ground-truth motion capture data (GT) as reference. Notably, our method can correct artifacts present in motion capture data, such as ground penetration or floating. The success rates show that our method leads to most stable grasps in the physics simulation. * The ground distance (GD) in the SMPL-X space is not zero as a consequence of the rigid body approximation of the human in the physics simulation. This metric equates to 0.0 when evaluated directly in the physics engine.

5.2.1 Baselines

Our method is capable of modeling the entire task of approaching an object, grasping and manipulating it. In contrast, the relevant baselines focus on a particular phase, e.g., GOAL [46] generates motions for the *approaching* phase while IMoS [9] tackles object *manipulation* after grasping. Hence, we compare our method against one baseline from each phase for a fair comparison. Though related, [23] is a very recent submission with no code publicly available.

We evaluate the baselines using the publicly available source code, pre-trained models and following the proposed evaluation protocols. Please note that there are differences between the settings of our method and IMoS. We model the entire task with a focus on single-handed object manipulation by providing an explicit control on the target trajectories. On the other hand, IMoS introduces language based control for two-handed object manipulation. Despite these differences, we deem a comparison justified since the physics-based metrics we report are invariant to the setting.

5.2.2 Metrics

We use the metrics proposed in prior works [6, 19, 46, 49]. The formal definitions are provided in supp. material.

Grasp Success Rate: We consider a grasp a success when the object is held for at least 0.5s in the physics simulation without dropping. For our model this includes approaching the object and lifting it from the table. We determine the success rate of the kinematics baselines using a static pose as a reference in physics simulation. The humanoid character and object are initialized with the last generated motion frame and maintain the grasp via PD-control [6, 19].

Ground Distance (GD): We compute the distance between

the average floating height (above ground) and the average vertical ground penetration depth, which are determined by the lowest SMPL-X vertex.

Foot Skating (FS): The percentage of foot skating frames. We consider a foot to be skating if the lowest SMPL-X vertex exceeds a threshold velocity [46].

Interpenetration: We report the interpenetration volume (IV) of MANO vertices that penetrate the object mesh and the maximum interpenetration depth (ID). In the *approaching* phase, we average the metric across the last five frames to be able to capture interpenetration before reaching the final grasp. For the *manipulation* phase, we average over five evenly distributed frames.

Trajectory Targets Reached (TTR): The ratio of the targets reached over all the targets in the trajectory. If a target is not reached within a certain time window, it is considered a failure and the next target from the trajectory is sampled. This metric is only applicable to our method and in the *manipulation* phase.

Contact Ratio (CR): The ratio of hand vertices that are within 5mm of the object mesh averaged over the sequence.

5.3. Evaluation

We provide a qualitative results of our method in Fig. 3 and a comparison against the baselines in Fig. 4. Please see our supplementary video for more examples.

We compare our method with GOAL [46] in the approaching phase until grasping and with IMoS [9] in the manipulation phase after grasping. Note that while we evaluate each phase separately, our method always performs the full sequence. We report the results in Tab. 2 using the metrics outlined in Section 5.2.2. We also provide the metrics for the ground truth (GT) as reference.

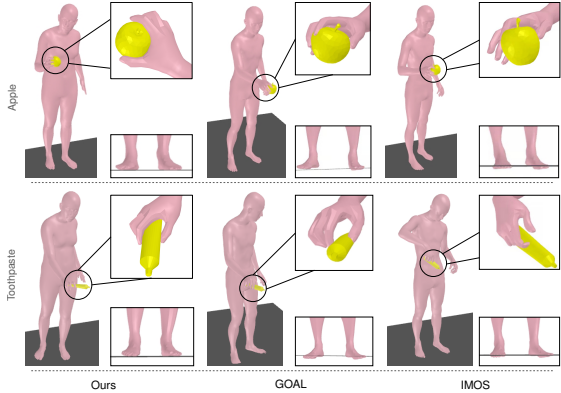


Figure 4. **Qualitative Comparison.** Our physics-based method generates motions that exhibit less hand-object interpenetration and ground interpenetration than the kinematics baselines.

Physical Plausibility Our method outperforms both baselines in all metrics, highlighting benefits of having a physics simulation in-the-loop. It leads to fewer artifacts as indicated by the hand-object interpenetration volume (IV) and depth (ID), foot skating (FS), and ground distance (GD). Baseline results often exhibit ground penetration, floating above ground, and hand-object collisions (see Fig. 4). Notably, our method also displays better physics-based properties compared to the ground truth data, which we argue is due to noise in the motion capture and labeling. Note that as a consequence of the approximated collision geometry as rigid bodies in the physics simulation, our method can still exhibit small amounts of interpenetration after converting the simulation results to the SMPL-X parameter space.

Contact Ratio To be in line with related work, we report the contact ratio (CR). We find that ours has a lower CR in the approaching phase than GOAL and a comparable CR with IMOS in the manipulation phase. However, we argue that this metric may not correlate with grasp quality due to the wide range of grasps. For example, grasps that mainly involve fingertips, such as a pinch grasp, lead to a lower CR. Furthermore, we observe that GOAL sometimes penetrates the object while approaching, yielding a high contact ratio despite the violation of physical constraints.

Success Rate Our method consistently achieves higher grasp success rates compared to the baselines. Note that simulation-based metrics such as grasp success have been established in previous works [6, 19] and give an indication on grasp stability. However, it should be interpreted with care when comparing physics and kinematic methods directly, since physics-based methods leverage a simulation, whereas kinematic-based methods do not. Small amounts of noise in contacts may already cause failure, because the PD-controller only maintains the input pose. Lastly, we validate how successful our method can follow a given target trajectory (TTR). The results indicate that most targets of the unseen test trajectories can be reached.

<i>decoupling</i>	<i>two-stage</i>	<i>t-guid.</i>	Success \uparrow	TTR \uparrow
\times	\times	\times	0.0	0.0
\checkmark	\times	\times	0.55	0.56
\checkmark	\checkmark	\times	0.77	0.79
\checkmark	\checkmark	\checkmark	0.79	0.85

Table 3. **Ablations.** We ablate the components of our method. The decoupling of priors is crucial to solve the task, while the two-stage training procedure and target guidance each contribute to higher success rates in grasping and trajectory following.

Generalization Our method can generalize to unseen objects (GOAL test set). It has difficulties grasping large objects where the fingers need to be fully stretched such as the *large cube* or *piggybank*. While these objects are part of the training set, they influence the success rate on the S10 test set. Examples of failure cases are in supp. material.

5.4. Ablations

We report ablation results in Tab. 3. We analyze the decoupling of the body prior from hand prior (*decoupling*), the two-stage training (*two-stage*) and the target guidance (*t-guid.*). We train all policies on the entire training set and evaluate on the test set. We find that decoupling of the coarse body motion from the dexterous hand motion is a critical component. Training a full-body prior directly leads to mode collapse in the latent space and hence fails to learn the full-body grasping task. The two-stage training procedure also plays an important role in achieving better performance. It allows the hand-object policy to first focus on achieving a stable grasp and then learn to follow the target trajectory. Lastly, our target guidance technique further improves the performance due to the explicit conditioning on target positions and the auxiliary training objective.

6. Discussion and Conclusion

We have introduced the first method to achieve physics-based full-body dexterous grasping. Our approach involves a hierarchical framework, beginning with the training of decoupled skill priors for body and hand control. These priors are then leveraged to develop a high-level policy to orchestrate the approaching, grasping and trajectory-guided manipulation phases. Notably, our method demonstrates a promising degree of physical plausibility in comparison to kinematics-based baselines. Our work also opens the door to potential future directions. For instance, there is potential in conditioning policies on language prompts, as shown in [9, 20], to guide the humanoid character. Moreover, our existing model relies on a single hand reference pose for guidance, a limitation that we hope could be addressed in future work. Lastly, while our current focus remains on single hand grasping, learning how to achieve physics-based bi-manual full-body grasping remains an open challenge.

References

- [1] Emre Aksan, Manuel Kaufmann, and Otmar Hilliges. Structured prediction helps 3d human motion modelling. In *ICCV*, 2019. First two authors contributed equally. 3
- [2] Emre Aksan, Manuel Kaufmann, Peng Cao, and Otmar Hilliges. A spatio-temporal transformer for 3d human motion prediction. 2021. 3
- [3] Kevin Bergamin, Simon Clavet, Daniel Holden, and James Richard Forbes. Drecon: data-driven responsive control of physics-based characters. *ACM Transactions On Graphics (TOG)*, 38(6):1–11, 2019. 2
- [4] Zhe Cao, Hang Gao, Karttikeya Mangalam, Qi-Zhi Cai, Minh Vo, and Jitendra Malik. Long-term human motion prediction with scene context. In *ECCV*, pages 387–404. Springer, 2020. 3
- [5] Yu-Wei Chao, Jimei Yang, Weifeng Chen, and Jia Deng. Learning to sit: Synthesizing human-chair interactions via hierarchical control. In *Proceedings of the AAAI Conference on Artificial Intelligence*, pages 5887–5895, 2021. 3
- [6] Sammy Christen, Muhammed Kocabas, Emre Aksan, Jemin Hwangbo, Jie Song, and Otmar Hilliges. D-grasp: Physically plausible dynamic grasp synthesis for hand-object interactions. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2022. 2, 3, 7, 8
- [7] Sammy Christen, Nina Schmid, and Otmar Hilliges. Generalizing skill embeddings across body shapes for physically simulated characters. *Embodied AI Workshop at Computer Vision and Pattern Recognition*, 2023. 5
- [8] Guillermo Garcia-Hernando, Edward Johns, and Tae-Kyun Kim. Physics-based dexterous manipulations with estimated hand poses and residual reinforcement learning. In *IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, pages 9561–9568. IEEE, 2020. 3
- [9] Anindita Ghosh, Rishabh Dabral, Vladislav Golyanik, Christian Theobalt, and Philipp Slusallek. Imos: Intent-driven full-body motion synthesis for human-object interactions. In *Eurographics*, 2023. 2, 3, 6, 7, 8
- [10] Partha Ghosh, Jie Song, Emre Aksan, and Otmar Hilliges. Learning human motion models for long-term predictions. pages 458–466. IEEE, 2017. 3
- [11] Kehong Gong, Bingbing Li, Jianfeng Zhang, Tao Wang, Jing Huang, Michael Bi Mi, Jiashi Feng, and Xinchao Wang. Posetriplet: Co-evolving 3d human pose estimation, imitation, and hallucination under self-supervision. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 11017–11027, 2022. 3
- [12] Patrick Grady, Chengcheng Tang, Christopher D. Twigg, Minh Vo, Samarth Brahmabhatt, and Charles C. Kemp. ContactOpt: Optimizing contact to improve grasps. In *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2021. 3
- [13] Mohamed Hassan, Duygu Ceylan, Ruben Villegas, Jun Saito, Jimei Yang, Yi Zhou, and Michael Black. Stochastic scene-aware motion prediction. In *ICCV*, 2021. 3
- [14] Mohamed Hassan, Yunrong Guo, Tingwu Wang, Michael Black, Sanja Fidler, and Xue Bin Peng. Synthesizing physical character-scene interactions. In *SIGGRAPH Conf. Track*, 2023. 2, 3
- [15] Jonathan Ho and Stefano Ermon. Generative adversarial imitation learning. In *Advances in Neural Information Processing Systems*. Curran Associates, Inc., 2016. 2
- [16] Daniel Holden, Jun Saito, Taku Komura, and Thomas Joyce. Learning motion manifolds with convolutional autoencoders. In *SIGGRAPH Asia 2015 Technical Briefs*, New York, NY, USA, 2015. Association for Computing Machinery. 3
- [17] Daniel Holden, Jun Saito, and Taku Komura. A deep learning framework for character motion synthesis and editing. *ACM Trans. Graph.*, 35(4), 2016. 3
- [18] Siyuan Huang, Zan Wang, Puhao Li, Baoxiong Jia, Tengyu Liu, Yixin Zhu, Wei Liang, and Song-Chun Zhu. Diffusion-based generation, optimization, and planning in 3d scenes, 2023. 3
- [19] Hanwen Jiang, Shaowei Liu, Jiashun Wang, and Xiaolong Wang. Hand-object contact consistency reasoning for human grasps generation. In *Proceedings of the International Conference on Computer Vision*, 2021. 7, 8
- [20] Jordan Juravsky, Yunrong Guo, Sanja Fidler, and Xue Bin Peng. Padl: Language-directed physics-based character control. In *SIGGRAPH Asia 2022 Conference Papers*, New York, NY, USA, 2022. Association for Computing Machinery. 2, 8
- [21] Diederik P. Kingma and Jimmy Ba. Adam: A method for stochastic optimization. In *International Conference on Learning Representations (ICLR)*, 2015. 6
- [22] Jiye Lee and Hanbyul Joo. Locomotion-action-manipulation: Synthesizing human-scene interactions in complex 3d environments, 2023. 3
- [23] Quanzhou Li, Jingbo Wang, Chen Change Loy, and Bo Dai. Task-oriented human-object interactions generation with implicit neural representations. *arXiv preprint arXiv:2303.13129*, 2023. 2, 3, 6, 7
- [24] Zhengyi Luo, Ryo Hachiuma, Ye Yuan, and Kris Kitani. Dynamics-regulated kinematic policy for egocentric pose estimation. *Advances in Neural Information Processing Systems*, 34:25019–25032, 2021. 3
- [25] Zhengyi Luo, Shun Iwase, Ye Yuan, and Kris Kitani. Embodied scene-aware human pose estimation. In *NeurIPS*, 2022. 2, 3
- [26] Viktor Makoviychuk, Lukasz Wawrzyniak, Yunrong Guo, Michelle Lu, Kier Storey, Miles Macklin, David Hoeller, Nikita Rudin, Arthur Allshire, Ankur Handa, et al. Isaac gym: High performance gpu-based physics simulation for robot learning. *arXiv preprint arXiv:2108.10470*, 2021. 6, 3
- [27] Priyanka Mandikal and Kristen Grauman. Learning dexterous grasping with object-centric visual affordances. In *IEEE International Conference on Robotics and Automation (ICRA)*, 2021. 3
- [28] Josh Merel, Saran Tunyasuvunakool, Arun Ahuja, Yuval Tassa, Leonard Hasenclever, Vu Pham, Tom Erez, Greg Wayne, and Nicolas Heess. Catch & carry: reusable neural controllers for vision-guided whole-body tasks. *ACM Transactions on Graphics (TOG)*, 39(4):39–1, 2020. 2

- [29] Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, et al. Pytorch: An imperative style, high-performance deep learning library. *Advances in neural information processing systems*, 2019. 6
- [30] Georgios Pavlakos, Vasileios Choutas, Nima Ghorbani, Timo Bolkart, Ahmed 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)*, pages 10975–10985, 2019. 3
- [31] Xue Bin Peng, Pieter Abbeel, Sergey Levine, and Michiel van de Panne. Deepmimic: Example-guided deep reinforcement learning of physics-based character skills. *ACM Trans. Graph.*, 37(4):143:1–143:14, 2018. 2
- [32] Xue Bin Peng, Michael Chang, Grace Zhang, Pieter Abbeel, and Sergey Levine. Mcp: Learning composable hierarchical control with multiplicative compositional policies. *Advances in Neural Information Processing Systems*, 32, 2019.
- [33] Xue Bin Peng, Ze Ma, Pieter Abbeel, Sergey Levine, and Angjoo Kanazawa. Amp: Adversarial motion priors for stylized physics-based character control. *ACM Trans. Graph.*, 40(4), 2021. 2
- [34] Xue Bin Peng, Yunrong Guo, Lina Halper, Sergey Levine, and Sanja Fidler. Ase: Large-scale reusable adversarial skill embeddings for physically simulated characters. *ACM Trans. Graph.*, 41(4), 2022. 2, 3, 4, 5, 6, 1
- [35] Yuxzhe Qin, Yueh-Hua Wu, Shaowei Liu, Hanwen Jiang, Ruihan Yang, Yang Fu, and Xiaolong Wang. Dexmv: Imitation learning for dexterous manipulation from human videos. *arXiv preprint arXiv:2108.05877*, 2021. 3
- [36] Aravind Rajeswaran*, Vikash Kumar*, Abhishek Gupta, Giulia Vezzani, John Schulman, Emanuel Todorov, and Sergey Levine. Learning Complex Dexterous Manipulation with Deep Reinforcement Learning and Demonstrations. In *Proceedings of Robotics: Science and Systems (RSS)*, 2018. 3
- [37] John Schulman, Sergey Levine, Pieter Abbeel, Michael Jordan, and Philipp Moritz. Trust region policy optimization. In *International conference on machine learning*, pages 1889–1897. PMLR, 2015. 2
- [38] John Schulman, Philipp Moritz, Sergey Levine, Michael Jordan, and Pieter Abbeel. High-dimensional continuous control using generalized advantage estimation. *arXiv preprint arXiv:1506.02438*, 2015.
- [39] John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, and Oleg Klimov. Proximal policy optimization algorithms. *arXiv preprint arXiv:1707.06347*, 2017. 2
- [40] Soshi Shimada, Vladislav Golyanik, Weipeng Xu, and Christian Theobalt. Physcap: Physically plausible monocular 3d motion capture in real time. *ACM TOG*, 39(6):1–16, 2020. 3
- [41] Soshi Shimada, Vladislav Golyanik, Weipeng Xu, Patrick Pérez, and Christian Theobalt. Neural monocular 3d human motion capture with physical awareness. *ACM TOG*, 40(4):1–15, 2021. 3
- [42] C. Karen Liu Sirui Chen, Albert Wu. Synthesize dexterous nonprehensile pregrasp for ungraspable objects. In *SIG-GRAPH '23: Special Interest Group on Computer Graphics and Interactive Techniques Conference*, 2023. 2
- [43] Sebastian Starke, Yiwei Zhao, Taku Komura, and Kazi Zaman. Local motion phases for learning multi-contact character movements. *ACM TOG*, 39(4):54–1, 2020. 3
- [44] RS Sutton and AG Barto. Introduction to reinforcement learning (1st ed.), 1998. 4, 6
- [45] Omid Taheri, Nima Ghorbani, Michael J. Black, and Dimitrios Tzionas. GRAB: A dataset of whole-body human grasping of objects. In *European Conference on Computer Vision (ECCV)*, 2020. 6, 3
- [46] Omid Taheri, Vasileios Choutas, Michael J. Black, and Dimitrios Tzionas. GOAL: Generating 4D whole-body motion for hand-object grasping. In *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2022. 2, 3, 6, 7, 4
- [47] Purva Tendulkar, Dídac Surís, and Carl Vondrick. Flex: Full-body grasping without full-body grasps. In *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2023. 3
- [48] Jingbo Wang, Sijie Yan, Bo Dai, and Dahua Lin. Scene-aware generative network for human motion synthesis. In *CVPR*, pages 12206–12215, 2021. 3
- [49] Yan Wu, Jiahao Wang, Yan Zhang, Siwei Zhang, Otmar Hilliges, Fisher Yu, and Siyu Tang. Saga: Stochastic whole-body grasping with contact. In *Proceedings of the European Conference on Computer Vision (ECCV)*, 2022. 2, 3, 6, 7
- [50] Kevin Xie, Tingwu Wang, Umar Iqbal, Yunrong Guo, Sanja Fidler, and Florian Shkurti. Physics-based human motion estimation and synthesis from videos. In *ICCV*, pages 11532–11541, 2021. 3
- [51] Ye Yuan and Kris Kitani. Residual force control for agile human behavior imitation and extended motion synthesis. *Advances in Neural Information Processing Systems*, 33: 21763–21774, 2020. 2
- [52] Ye Yuan, Shih-En Wei, Tomas Simon, Kris Kitani, and Jason Saragih. Simpo: Simulated character control for 3d human pose estimation. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 7159–7169, 2021. 3
- [53] Haotian Zhang, Ye Yuan, Viktor Makoviychuk, Yunrong Guo, Sanja Fidler, Xue Bin Peng, and Kayvon Fatahalian. Learning physically simulated tennis skills from broadcast videos. *ACM Trans. Graph.* 3
- [54] He Zhang, Yuting Ye, Takaaki Shiratori, and Taku Komura. Manipnet: neural manipulation synthesis with a hand-object spatial representation. *ACM Transactions on Graphics (ToG)*, 40(4):1–14, 2021. 2
- [55] Hui Zhang, Sammy Christen, Zicong Fan, Luocheng Zheng, Jemin Hwangbo, Jie Song, and Otmar Hilliges. Artigrasp: Physically plausible synthesis of bi-manual dexterous grasping and articulation, 2023. 5
- [56] Xiaohan Zhang, Bharat Lal Bhatnagar, Sebastian Starke, Vladimir Guzov, and Gerard Pons-Moll. COUCH: Towards controllable human-chair interactions. 2022. 3
- [57] Yi Zhou, Connelly Barnes, Jingwan Lu, Jimei Yang, and Hao Li. On the continuity of rotation representations in neural networks. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 5745–5753, 2019. 3

Physically Plausible Full-Body Hand-Object Interaction Synthesis

Supplementary Material

We provide this **manuscript** and a **video** as supplementary material. The table of contents below contains the structure of this document. Our code and models will be made publicly available upon publication.

Contents

A. Method Details	1
A.1. Discriminator Observations	1
A.2. Body-Prior Reward Function	1
A.3. Hand-object State Features	1
A.4. Task Reward Function	2
B. Implementation Details	3
B.1. Simulation environment	3
B.2. Preprocessing	3
B.3. Training Setup	3
C. Experimental Details	3
C.1. Metric Details	4
D. Additional Experiments	4
E. Limitations	5
F. Ethics Statement	5

A. Method Details

A.1. Discriminator Observations

Hand Prior The hand-prior discriminator features $\phi_h^D(\mathbf{s}) = (\boldsymbol{\theta}_h, \dot{\boldsymbol{\theta}}_h, \mathbf{x}_h^D)$ are equal to the hand-prior state features $\phi_h(\mathbf{s})$ with the exception that only wrist-relative 3D joint positions \mathbf{x}_h^D of fingertips (instead of all joints) are used. This design choice is motivated by [34], which uses a pruned version of the full state for the discriminator.

Body Prior The body-prior discriminator features are similar to the body-prior state features and defined as $\phi_b^D(\mathbf{s}) = (\boldsymbol{\theta}_b^D, \dot{\boldsymbol{\theta}}_b^D, \mathbf{x}_b^D, \mathbf{h}_b, \dot{\mathbf{t}}_b)$. The terms $\boldsymbol{\theta}_b^D$ and $\dot{\boldsymbol{\theta}}_b^D$ represent the local (parent-relative instead of root-relative as in $\phi_b(\mathbf{s})$) joint orientations and their angular velocities (except for the global root joint orientation and velocity). The root-relative 3D joint positions of key joints (left and right: elbow, wrist, knee, ankle, foot) are indicated by \mathbf{x}_b^D . The height of the root is defined by \mathbf{h}_b and the linear velocity of the 3D root position is given by $\dot{\mathbf{t}}_b$.

A.2. Body-Prior Reward Function

Besides the discriminator and encoder rewards outlined in Eq. (2) of the main paper, the body prior uses a trajectory reward r_ξ^b and a regularization reward r_{reg}^b :

$$r_b = r_\xi^b + r_{\text{reg}}^b - \log(1 - D(\phi_b(\mathbf{s}), \phi_b(\mathbf{s}'))) + \beta \log q(\mathbf{z}_b | \phi_b(\mathbf{s}), \phi_b(\mathbf{s}')). \quad (10)$$

Trajectory Reward Given a randomly sampled 3D target root position \mathbf{t}_b^i and target wrist position \mathbf{t}_h^i , the trajectory reward for the body prior is computed as the distance to the current root position \mathbf{t}_b and wrist position \mathbf{t}_h :

$$r_\xi^b = \exp(-\alpha_b \Delta t_b) + \exp(-\alpha_h \Delta t_h), \quad (11)$$

$$\Delta t_b = \left(\min \left(\|\mathbf{t}_b^i - \mathbf{t}_b\|, \beta_b \right) - \beta_b \right), \quad (12)$$

$$\Delta t_h = \left(\min \left(\|\mathbf{t}_h^i - \mathbf{t}_h\|, \beta_h \right) - \beta_h \right), \quad (13)$$

where the weights are defined by $\beta_b = 0.2$, $\beta_h = 0.005$, $\alpha_b = 2.0$, and $\alpha_h = 3.0$.

Regularization Reward To prevent fast, unnatural movements we regularize the linear wrist velocity $\dot{\mathbf{t}}_h$:

$$r_{\text{reg}}^b = 1 - \min(\exp(\|\dot{\mathbf{t}}_h\| - 0.8), 10). \quad (14)$$

A.3. Hand-object State Features

Contact Features We now explain in more detail the contact features \mathbf{g}_c from Eq. (4) of the main paper:

$$\mathbf{g}_c = (\bar{\mathbf{c}}_h, \bar{\mathbf{c}}_h \ominus \mathbf{c}_h). \quad (15)$$

The first term $\bar{\mathbf{c}}_h \in \mathbb{R}^{16 \times 1}$ is a binary target contact vector, which indicates which hand joints (16 in total) should be in contact with the object according to the hand pose reference Ψ . The second term is a distance vector with binary values showing whether a target contact is achieved or not:

$$\bar{\mathbf{c}}_h \ominus \mathbf{c}_h = \mathbb{I}_{\bar{\mathbf{c}}_h = \mathbf{c}_h = 1}. \quad (16)$$

For each contact body in \mathbf{c}_h , the vector is 0 unless a target contact $\bar{\mathbf{c}}_h$ is achieved, in which case it is 1.

Motion-Phase The term $\boldsymbol{\eta} \in [0, 1]$ indicates which phase of the task the human character is in. To this end, we define a set of six discrete states using the following heuristics:

1. The distance between the wrist and object is above 0.5m.

2. The distance between the wrist and object is below 0.5m, but above 0.2m.
3. The distance between the wrist and object is below 0.2m.
4. The hand is in contact with the object.
5. The object is lifted from the table.
6. The vertical distance between the initial object position and the current position is larger than 3cm.

To encode these states into the phase variable, we simply quantize the interval and assign it to the states in increasing order (i.e., the first state is assigned 0.0, the second state 0.2, etc.).

A.4. Task Reward Function

The task reward r_T of the hand-object interaction policy (see Eq. (7) in the main paper) is a linear combination between the static grasp reward (Appendix A.4.1), the trajectory reward (Appendix A.4.2), and a regularization reward (Appendix A.4.3).

A.4.1 Static Grasp Reward

The static grasp reward incentivizes the policy to grasp the object firmly such that it does not slip out of the hand. The reward is split into joint position reward r_x , joint orientation reward r_θ , and a contact reward r_c .

Position Reward The position reward promotes moving the wrist and finger joints (including the fingertips) to the 3D target joint positions given by the hand pose reference Ψ . To make the 3D target joint positions invariant with respect to the object pose, we convert all joint positions into object-relative frame. Given the current 3D target joint positions $\bar{\mathbf{x}}_h^j$ and the current 3D target joint positions \mathbf{x}_h^j of each joint j , we compute:

$$\Delta \mathbf{x}_{a:b} = \sum_{j=a}^b \left(\min \left(\left\| \bar{\mathbf{x}}_h^j - \mathbf{x}_h^j \right\|, \beta_x \right) - \beta_x \right), \quad (17)$$

$$r_x = 0.5 \exp(-1.25 \Delta \mathbf{x}_{1:J}) + \exp(-1.5 \Delta \mathbf{x}_{1:1}), \quad (18)$$

where J is the total number of joints, $\beta_x = 0.01\text{m}$ is a constant and $j = 1$ indicates the wrist joint.

Orientation Reward The orientation reward r_θ incentivizes the policy to move the wrist and finger joints into the target orientations given by the hand pose reference Ψ . We make use of the geodesic norm to compute the reward. Given the current joint rotation \mathbf{q}_h^j and the target joint rotation $\bar{\mathbf{q}}_h^j$ as quaternion of each joint j (which we convert

from θ_h and $\bar{\theta}_h$), we compute:

$$\begin{aligned} \Delta q_h^j &= \arccos(2 (\mathbf{q}_h^j \circ \bar{\mathbf{q}}_h^j)^2 - 1) \\ \Delta q_h &= \min \left(\frac{1}{J} \sum_{j=1}^J \Delta q_h^j, \beta_\theta \right) - \beta_\theta \\ r_\theta &= \exp(-2.5 \Delta q_h), \end{aligned} \quad (19)$$

where \circ indicates quaternion multiplication, J is the total number of joints, $\beta_\theta = 0.1\text{rad}$ is a constant and $j = 1$ indicates the wrist joint.

Contact reward The contact reward r_c comprises three components: the *contact-mask* reward $r_{c,\text{mask}}$, the *contact force* reward $r_{c,\text{force}}$, and the *no-table-contact* reward $r_{c,\text{tab}}$:

$$r_c = r_{c,\text{mask}} + r_{c,\text{force}} + r_{c,\text{tab}} \quad (20)$$

The *contact-mask* reward guides the hand parts towards reaching the target contacts extracted from the hand pose reference Ψ :

$$c_{\text{mask}} = \left(\frac{\bar{\mathbf{c}}_h^\top \mathbf{c}_h}{\bar{\mathbf{c}}_h^\top \bar{\mathbf{c}}_h} + \frac{\bar{\mathbf{c}}_h^\top \mathbf{c}_h^{t-1}}{\bar{\mathbf{c}}_h^\top \bar{\mathbf{c}}_h} \right) \quad (21)$$

$$r_{c,\text{mask}} = 1 - \exp(-1.6 c_{\text{mask}}). \quad (22)$$

The term $\frac{\bar{\mathbf{c}}_h^\top \mathbf{c}_h}{\bar{\mathbf{c}}_h^\top \bar{\mathbf{c}}_h}$ computes the ratio of number of bodies in contact with the object according to the hand pose reference. \mathbf{c}_h^{t-1} is the binary contact vector from the previous physics simulation state. Hence, the second term in Eq. (21) promotes coherent contacts over time. An entry in \mathbf{c}_h is 1 if the net contact force for that joint body is larger than zero.

The *contact force* reward incentivizes the policy to apply enough force between the hand and the object to grasp it stably:

$$\eta_{\text{force}} = \min(\|\mathbf{f}_h\|, \beta_c m_o) \quad (23)$$

$$r_{c,\text{force}} = \exp(0.25 \eta_{\text{force}}) \mathbb{I}_{\bar{\mathbf{c}}_h = \mathbf{c}_h = 1}, \quad (24)$$

$$(25)$$

where m_o is the object's weight and $\beta_c = 15$ is a constant. In essence, the term η_{force} promotes forces being applied up to an empirically defined maximum net force. The reward is only added for joints that are supposed to be in contact according to the target contacts $\bar{\mathbf{c}}_h$, which is indicated by $\mathbb{I}_{\bar{\mathbf{c}}_h = \mathbf{c}_h = 1}$.

The *no-table-contact* reward promotes being in contact with the object while avoiding forces applied to the table:

$$r_{c,\text{tab}} = \mathbb{I}_{\mathbf{f}_h > 0} \mathbb{I}_{\mathbf{f}_{\text{tab}} = 0}, \quad (26)$$

where $\mathbb{I}_{\mathbf{f}_h > 0}$ is an indicator for hand-object contact and $\mathbb{I}_{\mathbf{f}_{\text{tab}} = 0}$ is an indicator that is 1 if there is no force applied to the table by neither the object nor the hand. Note that the reward is non-zero only if both conditions are true.

A.4.2 Trajectory Reward

Given the current 3D root position \mathbf{t}_b , the current root-relative wrist position \mathbf{t}_h , and the current i -th trajectory target positions $(\mathbf{t}_b^i, \mathbf{t}_h^i)$, we compute the reward as described in Eq. (11), but with different weights and an additional component:

$$r_\xi = \exp(-\alpha_b \Delta t_b) + \exp(-\alpha_h \Delta t_h) + \alpha_s N_{\text{success}}, \quad (27)$$

where $\beta_b = 0.01, \beta_h = 0.01, \alpha_b = 1.25, \alpha_h = 3.0, \alpha_s = 0.008$. The last term is used to counterbalance a drop in the position reward as soon as a target is reached and a subsequent target is sampled, because this may make the policy not pursue any targets. This reward term increases with the number of achieved targets N_{success} .

A.4.3 Regularization Reward

The regularization reward $r_{\xi, \text{reg}}$ is defined as follows:

$$r_{\text{reg}} = \exp(-\|\dot{\mathbf{t}}_o\|) + \exp(-\|\ddot{\mathbf{t}}_h\|). \quad (28)$$

We regularize the object’s linear velocity $\dot{\mathbf{t}}_o$ and the jerk of the hand $\ddot{\mathbf{t}}_h$ (computed with finite differences from \mathbf{t}_h).

B. Implementation Details

B.1. Simulation environment

The physics simulation environment contains the humanoid, the object and a table. We model the table as a floating box and the object using its mesh. The provided meshes in GRAB have a high vertex count. In order to reduce the computational complexity of collision detection, we decimate all meshes. We compute the object weight based on the mesh volume and a constant density. We base our humanoid on the neutral SMPL-X [30] human body model but exclude eyeballs and jaw. The skeleton of the humanoid is created by extracting the joint positions and kinematic tree of the SMPL-X body model. We add an actuator to each joint and limit the joints based on the distribution of the GRAB dataset [45]. Similar to [52], we create a rigid body mesh for every joint of the SMPL-X body model. The body meshes are built by assigning each vertex to the joint with the largest linear blend skinning weight and then computing a convex hull per joint. The weight of each body is computed using the volume of the mesh and a constant density. To simplify the computational complexity, we approximate the collision geometries of the rigid body meshes with boxes, cylinders, and capsules, with the exception of the ankles and feet. Since we focus on right-hand grasping, we remove the left hand’s finger joints from the humanoid.

As Isaac Gym [26] does not yet allow to determine the origin of the net contact force experienced by a rigid body,

we disable certain collisions in order to retrieve useful contact observations. All collision between the humanoid and table are disabled. Moreover, all self-collisions between hand joints are disabled during the training of the hand-object interaction policy. However, self-collisions of the fingers are enabled during pre-training of the hand prior, which should prevent learning skills that cause self-penetration.

B.2. Preprocessing

As our humanoid character in the physics simulation is based on the neutral SMPL-X model, we need to convert the subject specific GRAB data. We first align the feet with the ground by translating each frame of the motion by the distance of the lowest SMPL-X vertex to the ground, i.e., we either lift or lower the character. To align the object with the hand, we translate the object and table by the distance between the thumb joints of the subject-specific and the neutral characters’ motions. We determine the hand pose reference using a heuristic, where we choose the frame within a time-window after the initial hand-object contact with the highest number of hand-object contacts. To add variety to training, we add multiple hand pose references close to the chosen frame in time. Finally, we optimize the hand poses of the references using ContactOpt [12]. To generate target trajectories, we extract a set of wrist and root position targets that are 1/15s apart from the motion capture reference motions, starting from the initial frame of hand-object contact. In our experiments, we limit the reference motions to a length of 4s. Instead of using one single set of targets per trajectory during training, we shift a window over the motion clip, which yields multiple sets of targets.

B.3. Training Setup

We use a single 80GB A100 to train the body and hand prior and a 24GB RTX 3090 TI NVIDIA GPU to train the hand-object interaction policy. We simulate 8192 parallel environments when training the priors and 2048 parallel environments for the hand-object interaction policy. The policies are updated after sampling 32 steps in each environment, yielding batches of ~262k and ~65k samples for the priors and the hand-object interaction policy, respectively. We train the priors for 40k and the hand-object interaction policy for 190k epochs, which amounts to roughly 6 days and 7 days of training, respectively.

C. Experimental Details

Randomization We randomly sample hand pose references and target trajectories during training. To increase robustness, we add uniform noise of $[-30, 30]$ mm to the hand pose references and $[-2, 2]$ mm to the trajectory targets, respectively.

Method	GD [mm] (\downarrow)	FS [%](\downarrow)	IV [cm^3] (\downarrow)	ID [mm] (\downarrow)	CR
Approaching					
Ours w/o smoothing	2.2	2.2	0.15	1.9	0.035
Ours	2.1	5.7	0.11	1.1	0.026
Manipulation					
Ours w/o smoothing (S10 test set)	2.0	4.2	0.13	1.6	0.030
Ours (S10 test set)	1.8	0.9	0.22	2.7	0.053
Ours w/o smoothing (GOAL test set)	2.1	4.4	0.13	1.4	0.033
Ours (GOAL test set)	1.9	1.2	0.18	2.9	0.055

Table 4. **Evaluation.** We evaluate our model without Gaussian smoothing (*w/o smoothing*) and compare with the results from the main paper. Note that the success rate and trajectory targets reached (TTR) are not affected as they are evaluated in the physics simulation.

Object	Success \uparrow	TTR \uparrow
apple	0.95	0.91
binoculars	0.54	0.83
camera	0.89	0.85
mug	0.64	0.74
toothpaste	0.94	0.94

Table 5. **Success And Trajectory Imitation Evaluation.** We evaluate the average success rate and the ratio of trajectory targets reached (TTR) of our method on the GOAL test set. The objects and trajectories are unseen during training.

C.1. Metric Details

Grasp Success Rate: We consider an object grasp successful if the object does not drop to the ground or table within a time window of 0.5s. For the baselines, we directly initialize the sequences in the predicted grasping pose without a table and consider a grasp successful if the object does not drop to the ground within 0.5s.

Ground Distance (GD): Given the set of SMPL-X 3D vertices \mathcal{V}_i per frame i , we extract the z-coordinate of the lowest vertex as $z_i = \min_z(\mathcal{V}_i)$. We compute the metric as follows:

$$GD = \frac{\sum_i z_i \mathbb{I}_{z_i > 0}}{\sum_i \mathbb{I}_{z_i > 0}} - \frac{\sum_i z_i \mathbb{I}_{z_i < 0}}{\sum_i \mathbb{I}_{z_i < 0}}. \quad (29)$$

It computes the distance between the average floating height and the average ground penetration depth. If $\sum_i \mathbb{I}_{z_i > 0} = 0$ or $\sum_i \mathbb{I}_{z_i < 0} = 0$, we use 0 for that term.

Foot Skating (FS): Given the set of SMPL-X 3D vertices \mathcal{V}_i per frame i , we find the vertex with the lowest z-coordinate $\mathbf{v}_j = \operatorname{argmin}_z(\mathcal{V}_i)$. The foot is considered skating if the horizontal velocity $\dot{\mathbf{v}}_j > 1\text{cm}$ per frame as proposed in [46] (note that we ignore the z-component of $\dot{\mathbf{v}}_j$). We compute the percentage of frames that are foot skating

over all frames N_{tot} :

$$FS = \frac{\sum_i^{N_{\text{tot}}} \mathbb{I}_{\dot{\mathbf{v}}_j > 1\text{cm}}}{N_{\text{tot}}} \quad (30)$$

Interpenetration: The interpenetration volume (IV) is computed as the average volume of vertices \mathcal{V} penetrating the object mesh. The interpenetration depth (ID) is given by the maximum distance between penetrating vertices and the object surface. In the *approaching* phase, we average the metric across the last five frames to capture interpenetration before reaching the final grasp. For the *manipulation* phase, we average over five evenly distributed frames.

Trajectory Targets Reached (TTR): Let N_{tot} be the total count of all reached targets in the trajectory and N_{success} the number of targets that were reached within a given time horizon of 0.2s, then $TTR = N_{\text{success}}/N_{\text{tot}}$. We consider a target reached if the wrist position is within 12cm of the target.

Contact Ratio: The ratio of SMPL-X vertices \mathcal{V}_i per frame i that are within 5mm of the object mesh, averaged over the whole sequence.

D. Additional Experiments

We provide a more detailed evaluation of two experiments. First, we report the success rate and the trajectory targets reached (TTR) metrics per object of the test set. The results are shown in Tab. 5. We find that the unseen objects with the most complex shapes, binoculars and mug, have the lowest success rates with 0.54 and 0.64, respectively. A better representation of the object shapes may alleviate such issues in the future. Furthermore, we report the metrics without applying Gaussian smoothing to our method in Tab. 4 (*w/o smoothing*). We find that it helps to improve the ground distance (GD) metric in both the approaching and the manipulation phase. In the approaching phase, it shows less interpenetration. In the manipulation phase, foot skating is reduced when applying smoothing. Moreover, we find the qualitative results to be more visually appealing with smoothing.

E. Limitations

We extend the discussion about limitations of our work and potential future directions from Section 6. We consider a unified body shape in our work. Exploring how to vary body shapes is a relatively under-researched problem in physics-based character control and more research is required [7]. Moreover, we use decimation to approximate the object mesh and body shape in order to make the physics simulation sufficiently fast for training. This leads to small interpenetration when converting back to the SMPL-X parametric space. As physics simulations develop, training with more high-resolution meshes will also become feasible. Lastly, our policy struggles with large objects, where the hands have to be fully stretched to grasp. Creating a framework for physics-based two-handed grasping, such as [55], but for full-body characters may help to overcome such edge cases.

F. Ethics Statement

Our work is in the realm of generating realistic and natural human motion data in simulations. This has future implications in areas such as AR/VR, human-computer interaction (HCI), and robotics. Therefore, one has to be careful in the utilization of such data. While the protection of user data is not a direct concern, since the data we generate is purely synthetic, the downstream use of the data has to be carefully considered. For example, while the generated data may serve in the training of service robots for hospitals or elderly care, it may just as well be used to train military robots. Moreover, being able to generate realistic virtual motions could be misused for generating deep-fakes when combined with realistic rendering techniques. While we don't have direct control over the explicit use cases of our technology, we believe discussing potential misuses of the technologies are important. Furthermore, we hope that openly sharing this research, the code and its technical details contributes to understanding the technology and enable access to as many users as possible.