
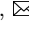


# AnySkill: Learning Open-Vocabulary Physical Skill for Interactive Agents

Jieming Cui<sup>1,2\*</sup>, Tengyu Liu<sup>2\*</sup>, Nian Liu<sup>2,3\*</sup>, Yaodong Yang<sup>1</sup>, Yixin Zhu<sup>1</sup>, , Siyuan Huang<sup>2</sup>, 

<sup>1</sup> Institute for Artificial Intelligence, Peking University    <sup>2</sup> National Key Laboratory of General Artificial Intelligence, BIGAI

<sup>3</sup> School of Artificial Intelligence, Beijing University of Posts and Telecommunications

<https://anyskill.github.io>

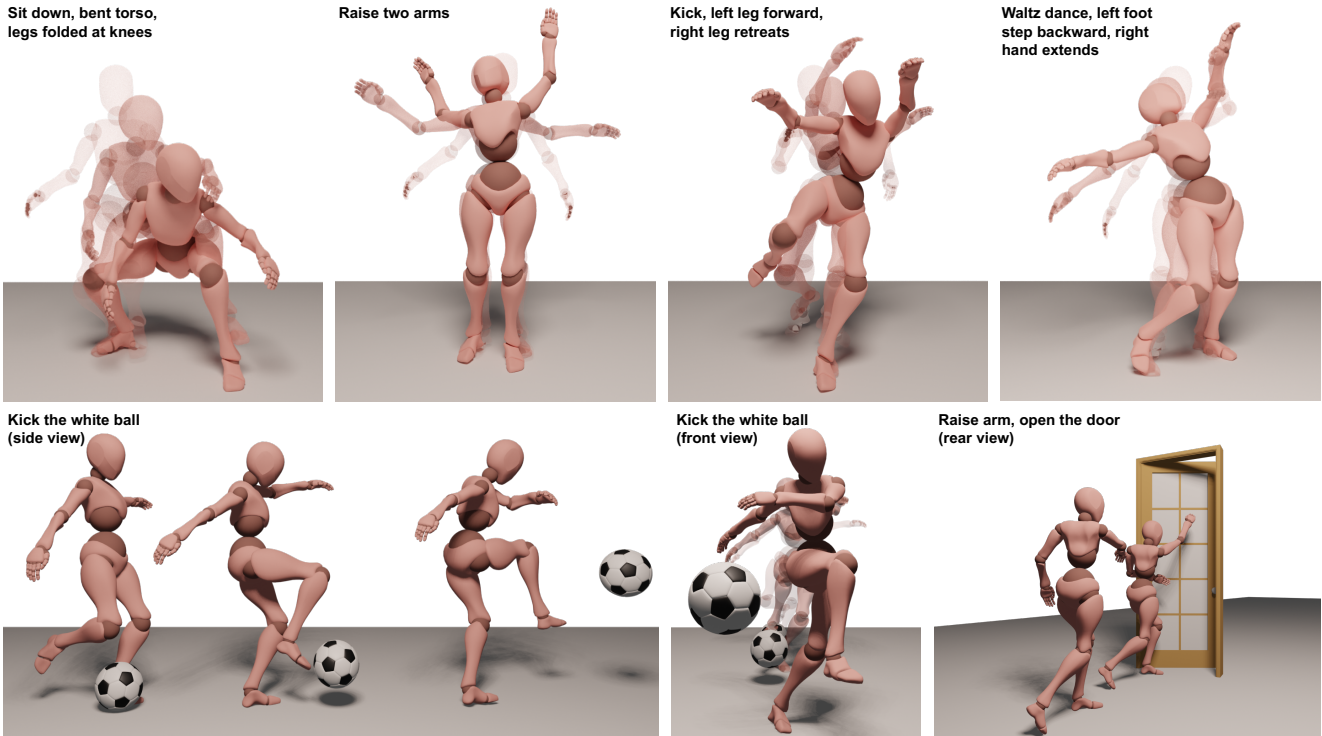


Figure 1. **Diverse motions generated by AnySkill conditioned on various instructions.** When provided with an open-vocabulary text description of a motion, AnySkill is adept at learning natural and flexible motions that closely align with the description, facilitated by an image-based reward mechanism. Additionally, AnySkill demonstrates proficiency in learning interactions with dynamic objects, showcasing its versatile motion generation capabilities.

## Abstract

Traditional approaches in physics-based motion generation, centered around imitation learning and reward shaping, often struggle to adapt to new scenarios. To tackle this limitation, we propose AnySkill, a novel hierarchical method that learns physically plausible interactions following open-vocabulary instructions. Our approach begins by developing a set of atomic actions via a low-level controller trained via imitation learning. Upon receiving an open-vocabulary textual instruction, AnySkill employs a high-level policy that selects and integrates these atomic actions to maximize the CLIP similarity between the agent’s rendered images and the text. An important feature of our method is the use of

image-based rewards for the high-level policy, which allows the agent to learn interactions with objects without manual reward engineering. We demonstrate AnySkill’s capability to generate realistic and natural motion sequences in response to unseen instructions of varying lengths, marking it the first method capable of open-vocabulary physical skill learning for interactive humanoid agents.

## 1. Introduction

Confronted with a soccer ball, an individual might engage in various actions such as kicking, dribbling, passing, or shooting. This interaction capability is feasible even for someone who has only observed soccer games, never having

played. This ability exemplifies the human aptitude for learning open-vocabulary physical interaction skills from visual experiences and applying these skills to novel objects and actions. Equipping interactive agents with this capability remains a significant challenge.

Recent physical skill learning methods predominantly rely on imitation learning to acquire realistic physical motions and interactions [29, 31]. However, this approach limits their adaptability to unforeseen scenarios with novel instructions and environments. Furthermore, neglecting physical laws in current models leads to unnatural and unrealistic motions, such as floating, penetration, and foot sliding, despite attempts to integrate physics-based penalties like gravity [58, 64] and collision [13, 57, 66]. Enhancing the generalizability of physically constrained motion generation is essential for decreasing reliance on specific datasets and fostering a more profound comprehension of the world.

On top of generalizability, the ultimate goal is to generate natural and interactive motions from any text input, known as achieving open vocabulary, which significantly increases the complexity of the problem. Several studies have explored open-vocabulary motion generation using large-scale pretrained models [11, 19, 37, 43]. However, these models struggle to produce natural motions, particularly interactive motions that require understanding broader environmental contexts or object interactions [11, 19, 43].

We identify a gap in motion generalizability on novel tasks and interaction capabilities with environments, hypothesizing that this is due to the reliance on improvised state representations and manually crafted reward mechanisms in prior works. Inspired by the human ability to learn new physical skills from visual inputs, we propose utilizing a Vision-Language Model (VLM) to offer flexible and generalizable state representations and image-based rewards for open-vocabulary skill learning. We introduce *AnySkill*, a hierarchical framework designed to equip virtual agents with the ability to learn open-vocabulary physical interaction skills. *AnySkill* combines a shared low-level controller with a high-level policy tailored to each instruction, learning a repertoire of latent atomic actions through generative adversarial imitation learning (GAIL), following CALM [42]. This ensures the naturalness and physical plausibility of each action. Then, for any open-vocabulary textual instruction, a high-level control policy dynamically selects latent atomic actions to optimize the CLIP [35] similarity between the agent’s rendered images and the textual instruction. This policy maintains physical plausibility and allows the agent to act according to a broad range of textual instructions. By leveraging CLIP similarity as a flexible and straightforward reward mechanism, our approach overcomes environmental limitations, facilitating interaction with any object. Despite the advances, creating natural and interactive actions for open-vocabulary models remains an ongoing challenge.

Extensive experiments demonstrate *AnySkill*’s ability to execute physical and interactive skills learned from open-vocabulary instructions; Fig. 1 showcases various interactive and non-interactive examples. We further prove that our method outperforms existing open-vocabulary motion generation approaches in creating interaction motions.

To summarize, our contributions are three-fold:

- We introduce *AnySkill*, a hierarchical approach that combines a low-level controller with a high-level policy, specifically designed for the learning of open-vocabulary physical skills.
- We leverage the VLM (*i.e.*, CLIP) to provide a novel means of generating flexible and generalizable image-based rewards. This approach eliminates the need for manually engineered rewards, facilitating the learning of both individual and interactive actions.
- Through extensive experimentation, we demonstrate that our method significantly surpasses existing approaches in both qualitative and quantitative measures. Importantly, *AnySkill* empowers agents with the ability to engage in smooth and natural interactions with dynamic objects across a variety of contexts.

## 2. Related Work

**Physical skills learning** emphasizes mastering motions that adhere to physical laws, including gravity, friction, and penetration. This domain has seen approaches that either employ specific loss functions to address constraints like foot-ground penetration [60], body-object interaction [1, 5, 8, 15, 21, 34, 47–50, 52, 59, 63, 65], self-collision [18, 27, 45], and gravity [6, 38, 54], or leverage physics simulators [16, 24, 31, 32, 42, 46] for more dynamic fidelity. Despite these efforts, ensuring fine-grained physical plausibility, especially in complex interactions, remains a challenge. The integration of reinforcement learning (RL) [10, 26, 29] and advanced modeling techniques (*e.g.*, MoE [2, 12, 53], VAE [20, 25], and GAN [10, 41]) alongside CLIP features [19, 37] attempts to improve generalization, yet faces the grand challenge of achieving physical plausibility in open vocabulary. Our method combines a shared low-level controller with a high-level policy tailored to each instruction, ensuring actions are physically realistic and adaptable to diverse instructions.

**Open-vocabulary motion generation** creates human motions from natural language descriptions outside the training distribution. Leveraging large-scale motion-language datasets [7, 23, 33], generative models have shown promise in motion synthesis [14, 36, 44, 62, 64]. However, these models often struggle with zero-shot generalization or adhering to the laws of physics, limited by their training data scope. Attempts to address these limitations include simplifying complex instructions with Large Language Models [17, 19] and employing pretrained VLMs like CLIP for

supervision [11, 22, 43], yet achieving natural and physics-compliant motions remains a significant hurdle. Our method builds upon these foundations, seeking to generate interactive and physically plausible motions from open-vocabulary descriptions, distinguishing itself from approaches like VLM-RMs [37] by modeling motion priors more effectively.

**Humanoid object interaction**, a relatively uncharted territory in physics-based motion generation, has seen simplifications such as attaching objects to characters’ hands to bypass the complexity of modeling physical interactions [29, 56, 61]. For dynamic interactions, encoding object states (positions and velocities) into the agent’s observations has facilitated specific tasks like dribbling [30, 31] and interacting with furniture [9], albeit requiring precise, object-specific rewards. This state-based approach is less feasible in open environments with diverse objects. Alternatively, vision-based policies [26] have shown potential for broader applications but are limited by their training domains. Our approach leverages a VLM for a more generalized motion-text alignment, avoiding the intricacies of manual reward crafting for varied interactive tasks.

### 3. AnySkill

AnySkill consists of two core components: the **low-level controller** and the **high-level policy**, illustrated in Fig. 2. Initially, we train a shared low-level controller,  $\pi^L$ , using unlabeled motion clips to distill a latent representation of atomic actions. This process utilizes GAIL [10], guaranteeing that the atomic actions are physically plausible.

Subsequently, for each open-vocabulary textual instruction, we train a high-level policy,  $\pi^H$ , tasked with composing atomic actions derived from low-level controllers. This high-level policy leverages a **flexible and generalizable image-based reward** via a VLM. This design facilitates the learning of physical interactions with dynamic objects, obviating the need for handcrafted reward engineering.

#### 3.1. Low-Level Controller

The low-level controller, inspired by CALM [42], enables the physically simulated humanoid agent to learn a diverse set of atomic actions. Formally, given an unlabeled motion dataset  $\mathcal{M}$ , we simultaneously train a motion encoder  $E$ , a discriminator  $D$ , and a controller  $\pi^L(a|s, z)$ . Here,  $a$  denotes the action,  $s$  the state, and  $z \in \mathcal{Z}$  the latent motion representation. The state  $s$  comprises the agent’s current root position, orientation, joint positions, and velocities, while the action  $a$  specifies the next target joint rotations.

Training proceeds as follows: A motion clip  $M$  from  $\mathcal{M}$  is encoded by  $E$  to yield the latent representation  $z = E(M)$ . The controller  $\pi^L(a|s, z)$  generates an action  $a$  based on the current state  $s$  and latent  $z$ . The agent then executes the action  $a$  in the physics-based simulator with a PD controller, resulting in a new state  $s'$ .

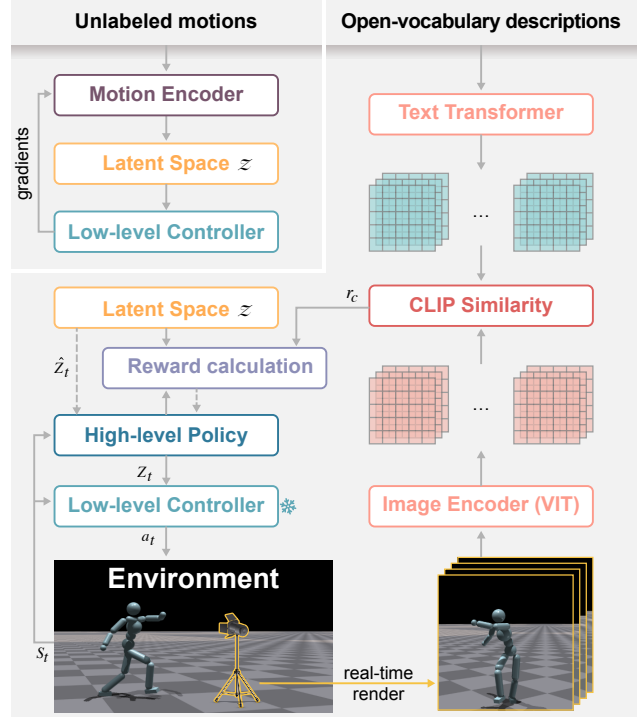


Figure 2. **The hierarchical structure of AnySkill.** Initially, the low-level controller (top-left) is trained to encode unlabeled motions into a shared latent space  $\mathcal{Z}$ . Subsequently, for each open-vocabulary text description, a high-level policy is trained. This policy orchestrates low-level actions to optimize the CLIP similarity between rendered images and the provided text, effectively composing actions that align with the textual instructions.

The discriminator  $D$  distinguishes whether the given  $(s, s')$  originates from the motion  $M$  corresponding to  $z$ , is produced by the controller  $\pi^L$  following the latent code  $z$ , or is produced by  $\pi^L$  following another latent code  $z' \sim \mathcal{Z}$ . We train  $D$  with a ternary adversarial loss:

$$\begin{aligned} \mathcal{L}_D = & -\mathbb{E}_{M \in \mathcal{M}} \left( \mathbb{E}_{d^\pi(s, s'|z)} [\log(1 - \mathcal{D}(s, s'|z))] \right) \quad (1) \\ & + \mathbb{E}_{d^M(s, s')} [\log \mathcal{D}(s, s'|z) + \log(1 - \mathcal{D}(s, s'|z' \sim \mathcal{Z}))] \\ & + w_{gp} \mathbb{E}_{d^M(s, s')} [ \|\nabla_{\theta} \mathcal{D}(\theta)|_{\theta=(s, s'|\hat{z})}\|^2 ] \Big|_{\hat{z} = \text{sg}(E(M))}, \end{aligned}$$

incorporating a gradient penalty with coefficient  $w_{gp}$  for stability, where  $\text{sg}(\cdot)$  denotes the stop gradient operator.

The encoder  $E$  is refined with both alignment and uniformity losses to ensure that embeddings of similar motions are closely aligned in the latent space, while dissimilar ones remain distinct [51], thus structuring  $\mathcal{Z}$  effectively.

The controller  $\pi^L$  aims to maximize the GAIL reward from  $D$ , calculated as

$$r^L(s, s', z) = -\log(1 - \mathcal{D}(s, s'|z)), \quad (2)$$

encouraging the generation of motions that closely resemble the original motion  $M$  associated with latent code  $z$ .

### 3.2. High-Level Policy

Building upon the atomic action repository created by the low-level controller, the high-level policy’s objective is to compose these actions, via the control of latent representation  $z$ , to generate motions that align with given text descriptions. With the low-level controller  $\pi^L$  fixed, we train a high-level policy  $\pi^H$  for each specific textual instruction, ensuring that the combined operation of both policy levels results in motions congruent with the text. The training process for the high-level policy is outlined in Algorithm 1.

---

#### Algorithm 1: Training of the high-level policy

---

**Input:** Reference motion dataset  $\mathcal{M}$ , frozen low-level controller  $\pi^L$ , frozen motion encoder  $E$ , simulation environment ENV, renderer image  $\mathcal{I}$ , CLIP feature of the description text  $f_d$

```

1  $\mathcal{Z} = E(\mathcal{M})$  initialize motion latent space
2 while not converged do
3    $\mathcal{B} \leftarrow \emptyset; p \leftarrow 0$  initialize
4   for horizon_length = 1, ...,  $n$  do
5     sample  $\hat{z}$  from  $\mathcal{Z}$ 
6     if horizon_length = 1 then
7        $s \leftarrow$  initialize;  $z \leftarrow \hat{z}$ 
8     else
9        $s \leftarrow$  ENV( $s, a$ );  $z \leftarrow \pi^H(s)$ 
10    end
11    for llc_steps = 1, ...,  $t$  do
12       $s \leftarrow$  ENV( $s, \pi^L(s, z)$ ) step simulation
13       $r^H \leftarrow$  calculate reward with Eq. (3)
14      if HEAD_HEIGHT < 0.15 then
15         $s, p \leftarrow 0$  reset agent and counter
16      end
17      if similarity is less than last step then
18         $p \leftarrow p + 1$  increment counter
19        if  $p \geq 8$  then
20           $p \leftarrow 0$  reset counter
21          reset  $s$  with 80% probability
22        end
23      end
24    end
25    update  $\mathcal{B}$  and  $\pi^H$  according to PPO
26  end
27 end

```

---

The high-level policy  $\pi^H$  is implemented as an MLP, taking the agent’s state  $s$  as input and outputting a latent representation  $z$  close to the low-level controller’s latent space  $\mathcal{Z}$ . It is trained using a composite reward of image-based similarity and latent-representation alignment. Given state  $s$  and text description  $d$ , we render the agent’s image  $\mathcal{I}(s)$  and encode it along with the text using a pretrained, frozen CLIP model to obtain features  $f_{\mathcal{I}}$  and  $f_d$ . The similarity reward is computed as the cosine similarity between  $f_{\mathcal{I}}$  and  $f_d$ , with an additional latent-representation alignment reward to

draw  $z$  nearer to the latent distribution of  $\mathcal{M}$ . The combined reward is given by:

$$r^H = \omega_c \cdot \frac{f_{\mathcal{I}} \cdot f_d}{|f_{\mathcal{I}}| |f_d|} + \omega_s \cdot \exp(-4\|z - \hat{z}\|_2), \quad (3)$$

where  $\omega_c, \omega_s$  are weighting factors, and  $\hat{z}$  is a sample from  $\mathcal{Z}$ . This image-based reward mechanism enables AnySkill to achieve text-to-motion alignment for open-vocabulary instructions. In addition, the image-based representation naturally encodes the entire environment around the agent, thus facilitating object interactions without modifying the encoding or architecture.

### 3.3. Implementation Details

**Low-level controller** The architecture of the encoder, low-level control policy, and discriminator comprises MLPs with hidden layers sized [1024, 1024, 512]. The latent space  $\mathcal{Z}$  is 64-dimensional. The alignment loss is set to 0.1, uniformity loss to 0.05, and gradient penalty to 5. The low-level controller is optimized using PPO [39] in IsaacGym. The training process is conducted on a single A100 GPU, operating at a 120Hz simulation frequency, and spans four days to cover a dataset comprising 93 unique motion patterns. Detailed hyperparameter settings of the low-level controller can be found in Tab. A1.

**High-level policy** The high-level policy, implemented as a two-layer MLP with hidden units of [1024, 512], outputs a 64-dimensional vector and is optimized using PPO. Training is conducted on an NVIDIA RTX3090 GPU, taking approximately 2.2 hours. Operationally, the high-level policy executes at a frequency of 6Hz, in contrast to the low-level policy, which operates at a more rapid 30Hz. This discrepancy in execution rates is strategic; the high-level policy is invoked every five timesteps, granting the low-level controller sufficient time to act on a given stable latent representation  $z$  and execute a complete atomic action. Such a setup is crucial for preventing the emergence of unnatural motion sequences by ensuring that each selected atomic action is fully realized before transitioning. Detailed hyperparameters of the high-level policy can be found in Tab. A2.

To further refine the training process and motion quality, an early termination strategy is employed to circumvent potential pitfalls of the high-level policy becoming trapped in suboptimal local minima. Specifically, the environment is reset with an 80% probability following eight successive reductions in CLIP similarity, or deterministically if the agent’s head height falls below 15cm. This approach significantly enhances training efficiency and the fidelity of the generated motions, ensuring a balance between exploration and the avoidance of poor performance traps.

**Rendering** We use IsaacGym’s default renderer, positioning the camera at (3m, 0m, 1m) while the agent is initialized at the origin. To maintain the agent at the focus

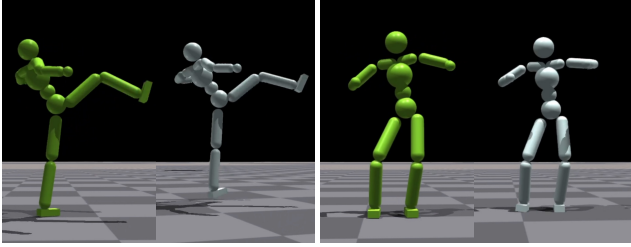


Figure 3. **Atomic actions from the trained low-level controller.** Each subfigure depicts the green agent demonstrating the reference motion from the dataset, while the white agent illustrates the corresponding learned atomic action.

of our visual feedback, we dynamically adjust the camera’s orientation each timestep to align with the agent’s pelvis joint. To encode the rendered images into a feature space compatible with our learning objectives, we employ the CLIP-ViT-B/32 model checkpoint from OpenCLIP [3], leveraging its robust representational capabilities.

**State projection** Given the computational demands of rendering images and extracting their CLIP features, we streamline the training process by introducing an MLP that projects the agent’s state vectors  $s$  directly to CLIP image features. This projection MLP is fine-tuned with an MSE loss against 104 million agent states accumulated during the high-level policy training. By substituting the render-and-encode steps with this MLP, we achieve a significant speedup, enhancing training efficiency by approximately 10.4 times, thereby mitigating the bottleneck associated with real-time image rendering and feature extraction.

## 4. Experiments

In this section, we detail the motion dataset curation for AnySkill’s low-level controller training (Sec. 4.1), evaluate AnySkill’s open-vocabulary motion generation against others (Sec. 4.2), analyze the text enhancement impact on effectiveness (Sec. 4.3), showcase physical interaction examples (Sec. 4.4), and compare our reward design with existing formulations (Sec. 4.5).

### 4.1. Training of Low-Level Controller

**Dataset** To enrich the low-level controller with diverse atomic actions, we assembled a dataset of 93 distinct motion records, primarily sourced from the CMU Graphics Lab Motion Capture Database [4] and SFU Motion Capture Database [40]. This collection spans various action categories, including locomotion (*e.g.*, walking, running, jumping), dance (*e.g.*, jazz, ballet), acrobatics (*e.g.*, roundhouse kicks), and interactive gestures (*e.g.*, pushing, greeting), all retargeted to a humanoid skeleton with 15 bones. We also adjusted any motions that lacked physical plausibility, ensuring the dataset’s fidelity for effective imitation learning.

**Training stabilization** Adversarial imitation learning’s instability, influenced by the volume and distribution of training data, can skew the density distribution in latent space, limiting the diversity of atomic actions for high-level policy selection. To mitigate this, we categorized motion records into 3 primary and 4 secondary groups by action scale and involved limbs. Details of the category division are described in Appendix A.2. By adjusting training data weights, we increased the likelihood of less frequent action groups, ensuring the variety of learned atomic actions; see also Fig. 3 and Fig. A7.

### 4.2. AnySkill Evaluation

Given the nascent field of open-vocabulary physical skill learning, we benchmark AnySkill against the two foremost similar methods in open-vocabulary motion generation: MotionCLIP [43] and AvatarCLIP [11], which also utilize CLIP similarity for generating human motions. To further understand the efficacy of our approach, we introduce a variant of our method, “Ours (no ET),” which operates without the early termination strategy.

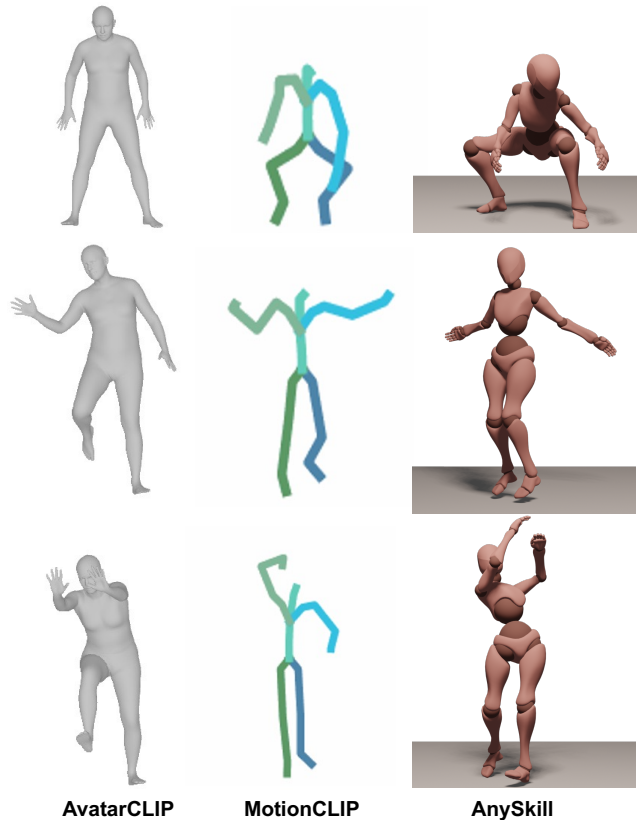


Figure 4. **Qualitative comparisons on open-vocabulary motion generation.** From top to bottom, the descriptions are “sit down, bent torso, legs folded at knees”, “legs off the ground, wave hands”, and “coiling the arm, throw a ball”. We showcase the most representative frames that best align with the descriptions.

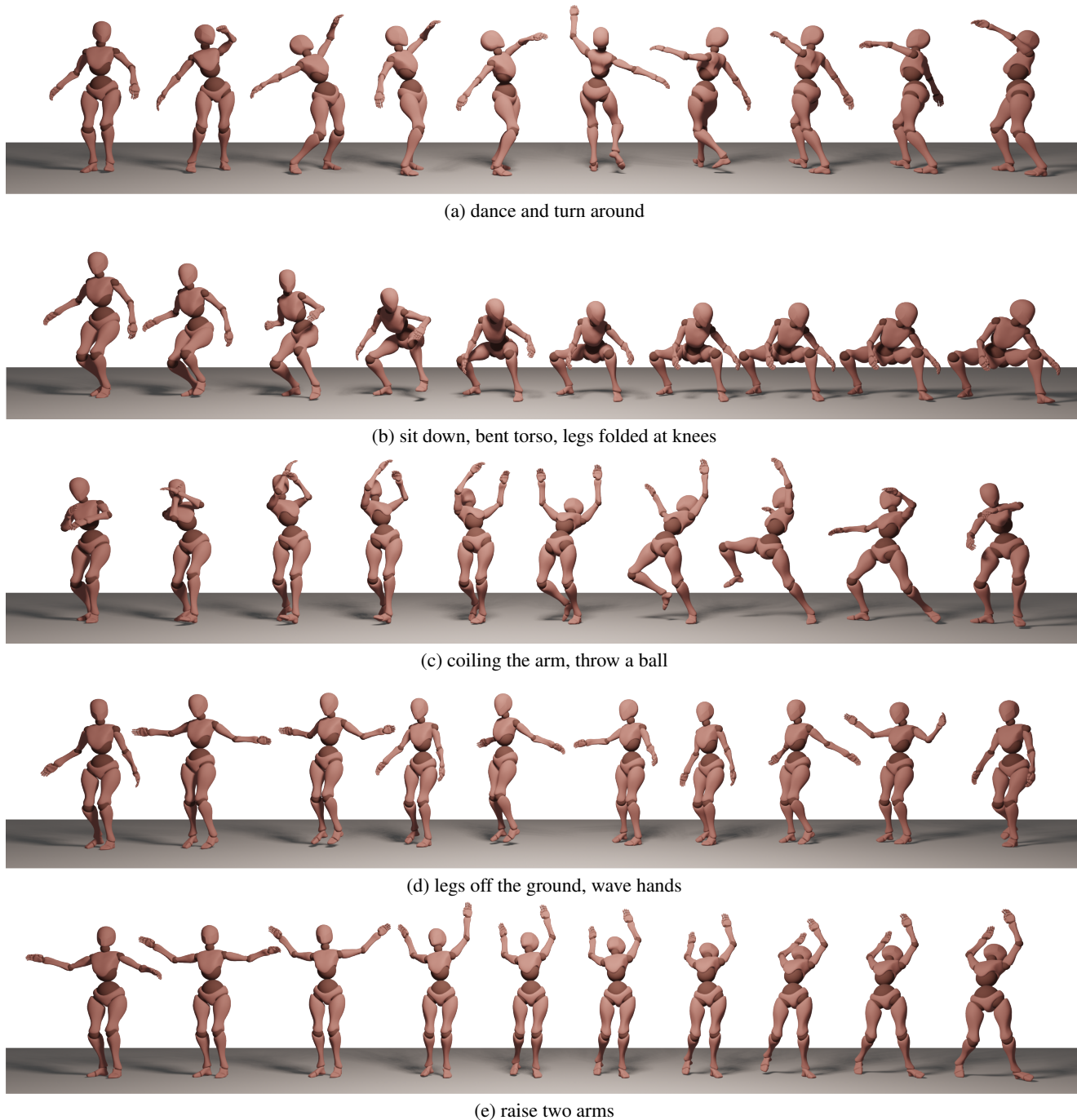


Figure 5. **Qualitative results of generated motion by AnySkill11.** Displayed are specific text descriptions and the corresponding motions generated by AnySkill11, as evaluated in the user study. Motion sequences progress from left to right.

For this evaluation, we selected 5 open-vocabulary text descriptions requiring comprehensive body movement and not covered in AnySkill’s training data. To assess the generated motions, we engaged 24 MTurk workers to rate them on task completion, smoothness, naturalness, and physical plausibility, using a scale from 0 to 10. Moreover, we computed the CLIP similarity score between the rendered

images and the text descriptions for each method as an objective measure. The motions generated by each method, including qualitative comparisons, are showcased in Fig. 4, with an in-depth look at AnySkill’s outputs presented in Fig. 5. Beyond the five actions presented, additional actions are shown in Appendix B.4

We present the results of the human study and quantita-

Table 1. Quantitative evaluation of high-level policy.

	Success ↑	Natural ↑	Smooth ↑	Physics ↑	CLIP_S ↑
AvatarCLIP [11]	4.29	4.74	5.79	5.74	21.11
MotionCLIP [43]	3.16	4.93	5.72	5.83	21.16
Ours (w/o ET)	5.05	4.88	5.68	5.31	21.89
Ours (w/o text-enhance)	3.06	4.48	5.19	5.96	20.76
Ours (w/ VideoCLIP [55])	2.37	4.90	5.65	6.41	21.35
<b>Ours (full)</b>	<b>6.16</b>	<b>6.23</b>	<b>6.51</b>	<b>6.93</b>	<b>24.18</b>

tive metrics in Tab. 1, demonstrating that AnySkill significantly surpasses current methods across all evaluated metrics. The ablation study underscores the importance of incorporating early termination into the training process. For additional comparative and qualitative results, see Fig. A5.

### 4.3. Text Enhancement

AnySkill excels at open-vocabulary skill acquisition, outperforming existing models. Its performance, however, is contingent on the specificity and scope of text descriptions. Performance drops with vague descriptions or for tasks requiring prolonged execution due to reliance on image-based similarity for rewards. For example, “do yoga” encompasses a broad range of poses, complicating convergence on a specific action. Similarly, for extended actions like “walk in a circle,” the model may not fully complete the task, as image-based rewards provide insufficient directional guidance.

To counteract these limitations, we introduced an automated script utilizing GPT-4 [28] to refine and clarify textual instructions, enhancing specificity and reducing potential motion interpretation ambiguity. This refinement process sig-

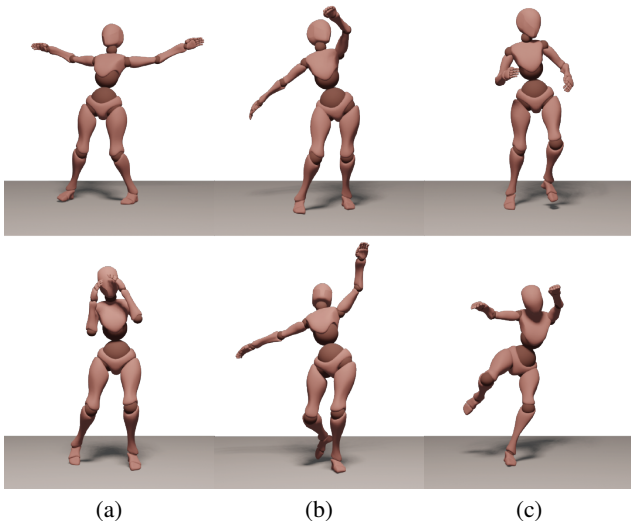


Figure 6. Qualitative evaluation of text description enhancement. We compare motions generated with original HumanML3D [7] descriptions (top row) against those from our enhanced descriptions (bottom row). Text descriptions are (a) “wave hi” and “raised arm bent at the elbow”; (b) “Waltz dance” and “left foot step backward, right hand extends”; (c) “kick” and “left leg forward, right leg retreats”.

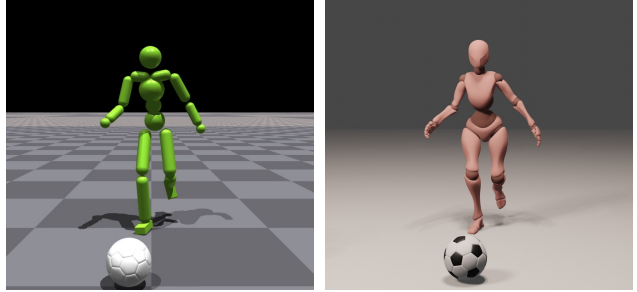


Figure 7. Agent and rendered mesh. The simulation of our agent and the interacting object (left) alongside their visualization (right).

nificantly improves AnySkill’s execution accuracy. Fig. 6 compares the original and refined texts alongside their generated motions; see Appendix B.1 for more qualitative results.

Moreover, we refined text descriptions from the HumanML3D [7] and BABEL [33] databases, amassing 1,896 unique, enhanced text instructions. For comprehensive details on the refined texts and their impact on motion generation, refer to Appendix A.1.

### 4.4. Interaction Motions

AnySkill demonstrates the superb capability to interact with dynamic objects, for instance, a soccer ball and a door. To capture these interactions accurately during training, we manually adjust the camera positions, focusing on the door and soccer ball. The alignment between the simulation environment and the rendered visualizations is showcased in Fig. 7. The qualitative assessments, as seen in Fig. 8, along with the quantitative evaluations in Tab. 2, confirm that AnySkill efficiently learns to interact with a variety of objects without necessitating any modifications to its learning algorithm or reward design. Our tests primarily involve interactions with a single object, yet extending AnySkill to engage with multiple objects concurrently is anticipated to be straightforward. Further interactive motions with various objects are available in Appendix B.4 and Fig. A8.

Table 2. Quantitative evaluation of interaction motions.

	Success ↑	Natural ↑	Smooth ↑	Physics ↑	CLIP_S ↑
Interaction w. object	5.42 -0.74	5.62 -0.61	5.34 -1.17	5.45 -1.48	24.49 +0.35
Interaction w. scene	4.53 -1.63	4.47 -1.76	5.01 -1.50	5.41 -1.52	22.41 -1.73

### 4.5. Reward Function Analysis

We evaluate 4 recent reward functions image- and physics-based RL and compare them with ours using cosine similarity. These include VLM-RMs [37], which adjusts the CLIP feature of text to exclude agent-specific details; CLIP-S [67], applying a modified CLIP similarity as the reward; VideoCLIP [55], calculating mean-pooled CLIP features across frames for temporal coherence; and ASE [32], adding a velocity reward for desired agent movement.

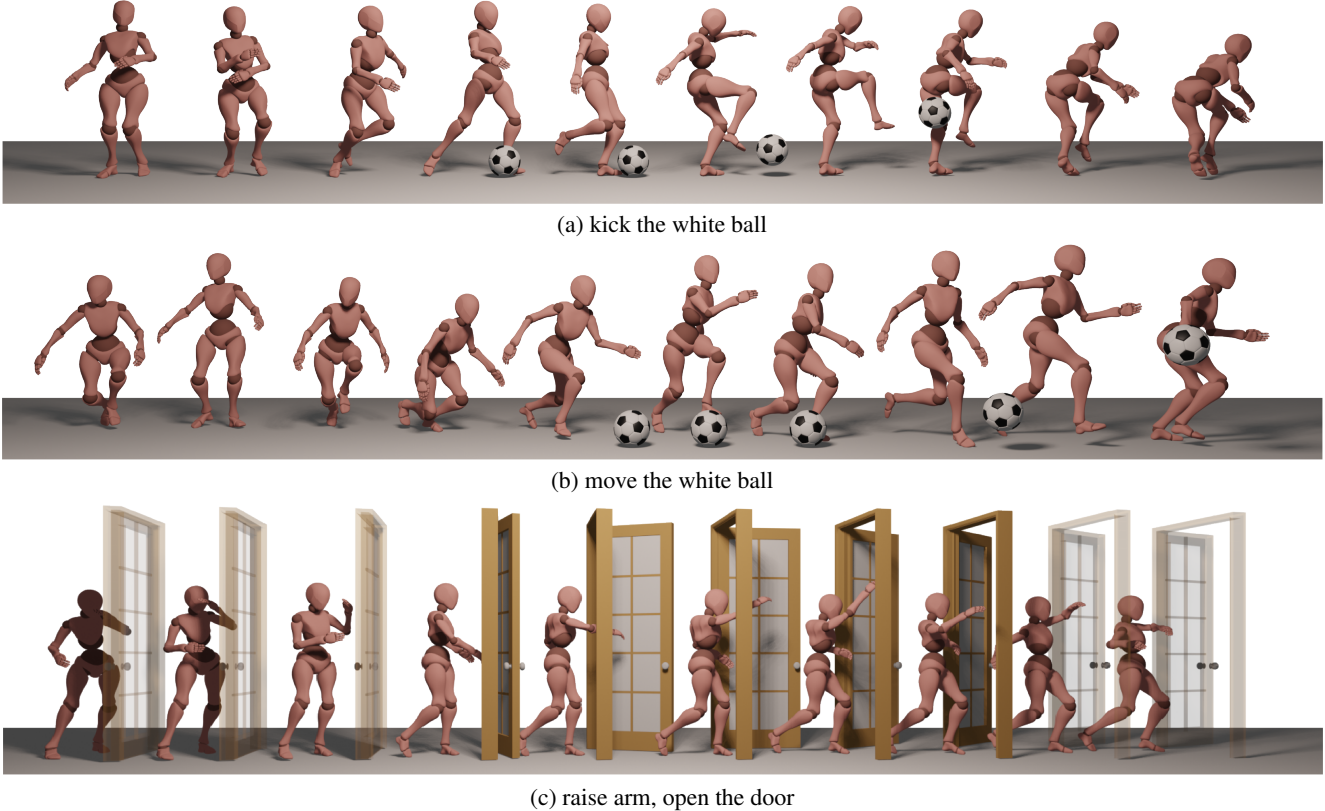


Figure 8. **Interaction motions generated by AnySkill.1.1.** Displayed are interaction sequences by AnySkill: two with a soccer ball (a-b) and one with a door (c), progressing from left to right.

Using these rewards, we train AnySkill on identical descriptions and assess motion quality via a user study similar to the one described in Sec. 4.2, with results presented in Tab. 3 and Appendix B.3. Our approach surpasses the baseline methods in most metrics, demonstrating the effectiveness of our reward function. Notably, AvgPool scores highly in smoothness, benefiting from averaging alignment scores over time.

Table 3. **Comparisons of the reward design.**

	Success $\uparrow$	Natural $\uparrow$	Smooth $\uparrow$	Physics $\uparrow$	CLIP_S $\uparrow$
VLM-RMs [37]	3.15	4.36	5.35	5.17	19.46
CLIP-S [67]	3.80	5.41	5.98	6.21	19.78
AvgPool [55]	5.09	5.96	<b>6.55</b>	6.70	20.25
+ vel. rew. [32]	2.73	4.42	5.35	5.22	18.39
<b>Ours</b>	<b>6.16</b>	<b>6.23</b>	6.51	<b>6.93</b>	<b>24.18</b>

## 5. Conclusion

We introduced AnySkill, a novel hierarchical framework for acquiring open-vocabulary physical interaction skills, combining an imitation-based low-level controller for motion generation with a robust, flexible image-based reward mechanism for adaptable skill learning. Through qualitative

and quantitative assessments, AnySkill is the first method capable of extending learning to encompass unseen tasks and interactions with novel objects, opening new venues in motion generation for interactive virtual agents.

**Future directions** AnySkill’s potential and limitations are closely linked to the CLIP model’s capabilities, guiding its current success and defining its challenges. As noted in Sec. 4.3, reliance on image-based rewards restricts AnySkill’s effectiveness in scenarios with prolonged durations or visual ambiguity. Future work aims to address these issues by enhancing the model’s understanding of temporal dynamics, integrating sophisticated multimodal alignment strategies, and incorporating interactive feedback loops.

The current need to develop a specialized policy for each new task—requiring substantial training time and resources—highlights a direction for future work: transforming AnySkill into a more universally applicable framework. This evolution will streamline the process of skill acquisition, dramatically reducing the time and resources required to master new interactive abilities. By achieving this, we anticipate enabling AnySkill to learn an array of skills in a unified, efficient manner, significantly broadening the scope of applications for interactive virtual agents and making sophisticated motion generation more accessible.



**Acknowledgement** The authors would like to thank Ms. Zhen Chen (BIGAI) for her exceptional contribution to the figure designs, Yanran Zhang and Jiale Yu (Tsinghua University) for their invaluable assistance in the experiments and prompt design, and Huiying Li (BIGAI) for crafting the agent’s appearance. We also thank NVIDIA for generously providing the necessary GPUs and hardware support. This work is supported in part by the National Science and Technology Major Project (2022ZD0114900), an NSFC fund (62376009), and the Beijing Nova Program.

## References

- [1] Yixin Chen, Siyuan Huang, Tao Yuan, Siyuan Qi, Yixin Zhu, and Song-Chun Zhu. Holistic++ scene understanding: Single-view 3d holistic scene parsing and human pose estimation with human-object interaction and physical commonsense. In *International Conference on Computer Vision (ICCV)*, 2019. 2
- [2] Zixiang Chen, Yihe Deng, Yue Wu, Quanquan Gu, and Yuanzhi Li. Towards understanding mixture of experts in deep learning. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2022. 2
- [3] Mehdi Cherti, Romain Beaumont, Ross Wightman, Mitchell Wortsman, Gabriel Ilharco, Cade Gordon, Christoph Schuhmann, Ludwig Schmidt, and Jenia Jitsev. Reproducible scaling laws for contrastive language-image learning. In *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2023. 5
- [4] CMU MOCAP. CMU Graphics Lab Motion Capture Database. <https://http://mocap.cs.cmu.edu/>, 2010. Accessed: 2023-10-25. 5
- [5] Jieming Cui, Ziren Gong, Baoxiong Jia, Siyuan Huang, Zilong Zheng, Jianzhu Ma, and Yixin Zhu. Probio: A protocol-guided multimodal dataset for molecular biology lab. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2024. 2
- [6] Erik Gärtner, Mykhaylo Andriluka, Erwin Coumans, and Cristian Sminchisescu. Differentiable dynamics for articulated 3d human motion reconstruction. In *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2022. 2
- [7] Chuan Guo, Shihao Zou, Xinxin Zuo, Sen Wang, Wei Ji, Xingyu Li, and Li Cheng. Generating diverse and natural 3d human motions from text. In *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2022. 2, 7, A1
- [8] Mohamed Hassan, Duygu Ceylan, Ruben Villegas, Jun Saito, Jimei Yang, Yi Zhou, and Michael J Black. Stochastic scene-aware motion prediction. In *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2021. 2
- [9] Mohamed Hassan, Yunrong Guo, Tingwu Wang, Michael Black, Sanja Fidler, and Xue Bin Peng. Synthesizing physical character-scene interactions. In *ACM SIGGRAPH Conference Proceedings*, 2023. 3
- [10] Jonathan Ho and Stefano Ermon. Generative adversarial imitation learning. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2016. 2, 3
- [11] Fangzhou Hong, Mingyuan Zhang, Liang Pan, Zhongang Cai, Lei Yang, and Ziwei Liu. Avatarclip: Zero-shot text-driven generation and animation of 3d avatars. *ACM Transactions on Graphics (TOG)*, 41(4):1–19, 2022. 2, 3, 5, 7, A2
- [12] Jiang Hua, Liangcai Zeng, Gongfa Li, and Zhaojie Ju. Learning for a robot: Deep reinforcement learning, imitation learning, transfer learning. *Sensors*, 21(4):1278, 2021. 2
- [13] Siyuan Huang, Zan Wang, Puhao Li, Baoxiong Jia, Tengyu Liu, Yixin Zhu, Wei Liang, and Song-Chun Zhu. Diffusion-based generation, optimization, and planning in 3d scenes. In *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2023. 2
- [14] Biao Jiang, Xin Chen, Wen Liu, Jingyi Yu, Gang Yu, and Tao Chen. Motiongpt: Human motion as a foreign language. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2024. 2
- [15] Nan Jiang, Tengyu Liu, Zhexuan Cao, Jieming Cui, Zhiyuan Zhang, Yixin Chen, He Wang, Yixin Zhu, and Siyuan Huang. Full-body articulated human-object interaction. In *International Conference on Computer Vision (ICCV)*, 2023. 2
- [16] Jordan Juravsky, Yunrong Guo, Sanja Fidler, and Xue Bin Peng. Padl: Language-directed physics-based character control. In *ACM SIGGRAPH Conference Proceedings*, 2022. 2
- [17] Sai Shashank Kalakonda, Shubh Maheshwari, and Ravi Kiran Sarvadevabhatla. Action-gpt: Leveraging large-scale language models for improved and generalized zero shot action generation. *arXiv preprint arXiv:2211.15603*, 2022. 2
- [18] Charles Khazoom, Daniel Gonzalez-Diaz, Yanran Ding, and Sangbae Kim. Humanoid self-collision avoidance using whole-body control with control barrier functions. In *International Conference on Humanoid Robots (Humanoids)*, 2022. 2
- [19] K Niranjan Kumar, Irfan Essa, and Sehoon Ha. Words into action: Learning diverse humanoid behaviors using language guided iterative motion refinement. In *2nd Workshop on Language and Robot Learning: Language as Grounding*, 2023. 2
- [20] Alex X Lee, Anusha Nagabandi, Pieter Abbeel, and Sergey Levine. Stochastic latent actor-critic: Deep reinforcement learning with a latent variable model. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2020. 2
- [21] Jiye Lee and Hanbyul Joo. Locomotion-action-manipulation: Synthesizing human-scene interactions in complex 3d environments. In *International Conference on Computer Vision (ICCV)*, 2023. 2
- [22] Junfan Lin, Jianlong Chang, Lingbo Liu, Guanbin Li, Liang Lin, Qi Tian, and Chang-wen Chen. Being comes from not-being: Open-vocabulary text-to-motion generation with wordless training. In *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2023. 3
- [23] Jing Lin, Ailing Zeng, Shunlin Lu, Yuanhao Cai, Ruimao Zhang, Haoqian Wang, and Lei Zhang. Motion-x: A large-scale 3d expressive whole-body human motion dataset. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2023. 2
- [24] 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. 2
- [25] Josh Merel, Leonard Hasenclever, Alexandre Galashov, Arun Ahuja, Vu Pham, Greg Wayne, Yee Whye Teh, and Nicolas Heess. Neural probabilistic motor primitives for humanoid control. In *International Conference on Learning Representations (ICLR)*, 2018. 2
- [26] 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, 3
- [27] Marko Mihajlovic, Shunsuke Saito, Aayush Bansal, Michael Zollhoefer, and Siyu Tang. Coap: Compositional articulated occupancy of people. In *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2022. 2
- [28] OpenAI. Introducing gpt-4. <https://openai.com/blog/gpt-4>, 2023. 7, A1
- [29] Xue Bin Peng, Pieter Abbeel, Sergey Levine, and Michiel Van de Panne. Deepmimic: Example-guided deep reinforcement learning of physics-based character skills. *ACM Transactions on Graphics (TOG)*, 37(4):1–14, 2018. 2, 3
- [30] Xue Bin Peng, Michael Chang, Grace Zhang, Pieter Abbeel, and Sergey Levine. Mep: Learning composable hierarchical control with multiplicative compositional policies. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2019. 3
- [31] Xue Bin Peng, Ze Ma, Pieter Abbeel, Sergey Levine, and Angjoo Kanazawa. Amp: Adversarial motion priors for stylized physics-based character control. *ACM Transactions on Graphics (TOG)*, 40(4):1–20, 2021. 2, 3
- [32] Xue Bin Peng, Yunrong Guo, Lina Halper, Sergey Levine, and Sanja Fidler. Ase: Large-scale reusable adversarial skill embeddings for physically simulated characters. *ACM Transactions on Graphics (TOG)*, 41(4):1–17, 2022. 2, 7, 8
- [33] Abhinanda R Punnakkal, Arjun Chandrasekaran, Nikos Athanasiou, Alejandra Quiros-Ramirez, and Michael J Black. Babel: Bodies, action and behavior with english labels. In *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2021. 2, 7, A1
- [34] Siyuan Qi, Yixin Zhu, Siyuan Huang, Chenfanfu Jiang, and Song-Chun Zhu. Human-centric indoor scene synthesis using stochastic grammar. In *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2018. 2
- [35] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual models from natural language supervision. In *International Conference on Machine Learning (ICML)*, 2021. 2
- [36] Zhiyuan Ren, Zhihong Pan, Xin Zhou, and Le Kang. Diffusion motion: Generate text-guided 3d human motion by diffusion model. In *IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, 2023. 2
- [37] Juan Rocamonde, Victoriano Montesinos, Elvis Nava, Ethan Perez, and David Lindner. Vision-language models are zero-shot reward models for reinforcement learning. In *International Conference on Learning Representations (ICLR)*, 2023. 2, 3, 7, 8
- [38] Tim Salzman, Marco Pavone, and Markus Ryll. Motron: Multimodal probabilistic human motion forecasting. In *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2022. 2
- [39] John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, and Oleg Klimov. Proximal policy optimization algorithms. *arXiv preprint arXiv:1707.06347*, 2017. 4
- [40] SFU MOCAP. SFU Motion Capture Database. <https://mocap.cs.sfu.ca/>, 2023. Accessed: 2023-11-01. 5
- [41] Jiaming Song, Hongyu Ren, Dorsa Sadigh, and Stefano Ermon. Multi-agent generative adversarial imitation learning. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2018. 2
- [42] Chen Tessler, Yoni Kasten, Yunrong Guo, Shie Mannor, Gal Chechik, and Xue Bin Peng. Calm: Conditional adversarial latent models for directable virtual characters. In *ACM SIGGRAPH Conference Proceedings*, 2023. 2, 3
- [43] Guy Tevet, Brian Gordon, Amir Hertz, Amit H Bermano, and Daniel Cohen-Or. Motionclip: Exposing human motion generation to clip space. In *European Conference on Computer Vision (ECCV)*, 2022. 2, 3, 5, 7, A2
- [44] Guy Tevet, Sigal Raab, Brian Gordon, Yonatan Shafir, Daniel Cohen-Or, and Amit H Bermano. Human motion diffusion model. In *International Conference on Learning Representations (ICLR)*, 2022. 2
- [45] Yating Tian, Hongwen Zhang, Yebin Liu, and Limin Wang. Recovering 3d human mesh from monocular images: A survey. *Transactions on Pattern Analysis and Machine Intelligence (TPAMI)*, 2023. 2
- [46] Emanuel Todorov, Tom Erez, and Yuval Tassa. Mujoco: A physics engine for model-based control. In *International Conference on Intelligent Robots and Systems (IROS)*, 2012. 2
- [47] Shashank Tripathi, Lea Müller, Chun-Hao P Huang, Omid Taheri, Michael J Black, and Dimitrios Tzionas. 3d human pose estimation via intuitive physics. In *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2023. 2
- [48] Jonathan Tseng, Rodrigo Castellon, and Karen Liu. Edge: Editable dance generation from music. In *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2023.
- [49] Jiashun Wang, Huazhe Xu, Jingwei Xu, Sifei Liu, and Xiao-long Wang. Synthesizing long-term 3d human motion and interaction in 3d scenes. In *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2021.
- [50] Jingbo Wang, Sijie Yan, Bo Dai, and Dahua Lin. Scene-aware generative network for human motion synthesis. In *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2021. 2
- [51] Tongzhou Wang and Phillip Isola. Understanding contrastive representation learning through alignment and uniformity on the hypersphere. In *International Conference on Machine Learning (ICML)*, 2020. 3
- [52] Zan Wang, Yixin Chen, Baoxiong Jia, Puhao Li, Jinlu Zhang, Jingze Zhang, Tengyu Liu, Yixin Zhu, Wei Liang, and Siyuan

- Huang. Move as you say, interact as you can: Language-guided human motion generation with scene affordance. In *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2024. 2
- [53] Jungdam Won, Deepak Gopinath, and Jessica Hodgins. A scalable approach to control diverse behaviors for physically simulated characters. *ACM Transactions on Graphics (TOG)*, 39(4):33–1, 2020. 2
- [54] Kevin Xie, Tingwu Wang, Umar Iqbal, Yunrong Guo, Sanja Fidler, and Florian Shkurti. Physics-based human motion estimation and synthesis from videos. In *International Conference on Computer Vision (ICCV)*, 2021. 2
- [55] Hu Xu, Gargi Ghosh, Po-Yao Huang, Dmytro Okhonko, Armen Aghajanyan, Florian Metze, Luke Zettlemoyer, and Christoph Feichtenhofer. Videoclip: Contrastive pre-training for zero-shot video-text understanding. *arXiv preprint arXiv:2109.14084*, 2021. 7, 8
- [56] Pei Xu, Xiumin Shang, Victor Zordan, and Ioannis Karamouzas. Composite motion learning with task control. *ACM Transactions on Graphics (TOG)*, 42(4):1–14, 2023. 3
- [57] Sirui Xu, Zhengyuan Li, Yu-Xiong Wang, and Liang-Yan Gui. Interdiff: Generating 3d human-object interactions with physics-informed diffusion. In *International Conference on Computer Vision (ICCV)*, 2023. 2
- [58] Shusheng Xu, Huaijie Wang, Jiaxuan Gao, Yutao Ouyang, Chao Yu, and Yi Wu. Language-guided generation of physically realistic robot motion and control. *arXiv preprint arXiv:2306.10518*, 2023. 2
- [59] Hongwei Yi, Chun-Hao P Huang, Shashank Tripathi, Lea Hering, Justus Thies, and Michael J Black. Mime: Human-aware 3d scene generation. In *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2023. 2
- [60] Ye Yuan, Jiaming Song, Umar Iqbal, Arash Vahdat, and Jan Kautz. Physdiff: Physics-guided human motion diffusion model. In *International Conference on Computer Vision (ICCV)*, 2023. 2
- [61] Haotian Zhang, Ye Yuan, Viktor Makovychuk, Yunrong Guo, Sanja Fidler, Xue Bin Peng, and Kayvon Fatahalian. Learning physically simulated tennis skills from broadcast videos. *ACM Transactions on Graphics (TOG)*, 42(4):1–14, 2023. 3
- [62] Jianrong Zhang, Yangsong Zhang, Xiaodong Cun, Shaoli Huang, Yong Zhang, Hongwei Zhao, Hongtao Lu, and Xi Shen. T2m-gpt: Generating human motion from textual descriptions with discrete representations. In *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2023. 2
- [63] Yan Zhang, Mohamed Hassan, Heiko Neumann, Michael J Black, and Siyu Tang. Generating 3d people in scenes without people. In *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2020. 2
- [64] Yaqi Zhang, Di Huang, Bin Liu, Shixiang Tang, Yan Lu, Lu Chen, Lei Bai, Qi Chu, Nenghai Yu, and Wanli Ouyang. Motiongpt: Finetuned llms are general-purpose motion generators. *arXiv preprint arXiv:2306.10900*, 2023. 2
- [65] Kaifeng Zhao, Shaofei Wang, Yan Zhang, Thabo Beeler, and Siyu Tang. Compositional human-scene interaction synthesis with semantic control. In *European Conference on Computer Vision (ECCV)*, 2022. 2
- [66] Kaifeng Zhao, Yan Zhang, Shaofei Wang, Thabo Beeler, and Siyu Tang. Synthesizing diverse human motions in 3d indoor scenes. In *International Conference on Computer Vision (ICCV)*, 2023. 2
- [67] Shuai Zhao, Xiaohan Wang, Linchao Zhu, and Yi Yang. Test-time adaptation with CLIP reward for zero-shot generalization in vision-language models. In *International Conference on Learning Representations (ICLR)*, 2024. 7, 8

## A. Data

This section offers a detailed account of the data’s origins and the methodologies employed for its processing.

### A.1. Text Data

Text descriptions sourced from publicly available online datasets are often marked by redundancy, ambiguity, and insufficient detail. To address these issues, it is necessary to preprocess the descriptions to render them more practical and usable. For generating practical text descriptions, we implemented a three-tiered process leveraging GPT-4 [28]. This encompasses **filtering text** to discard non-essential details, **scoring text** for assessing utility, and **rewriting text** to improve clarity and applicability. Our goal is to identify text descriptions that significantly contribute to mastering open-vocabulary physical skills from a robust pre-existing dataset, and to standardize the collection of text instructions.

**Filter text** Initially, we compiled 89,910 text entries from HumanML3D [7] and Babel [33], discovering substantial repetition, including exact duplicates, descriptions of akin actions (e.g., “A person walks down a set of stairs” vs. “A person walks down stairs”), frequency-related repetitions (e.g., “A person sways side to side multiple times” vs. “A person sways from side to side”), and semantic duplicates (e.g., “The person is doing a waltz dance” vs. “A man waltzes backward in a circle”).

To address this issue, we initiated a deduplication process, first eliminating descriptions that were overly brief (under three tokens) or excessively lengthy (over 77 tokens). We then utilized the LLAMA-2-7B MODEL with its 4096-dimensional embedding vector for further deduplication. By computing cosine similarities between each description pair and applying a 0.92 similarity threshold, descriptions exceeding this threshold were considered repetition. This procedure refined our dataset to 4,910 unique descriptions.

**Scoring text** After filtering out duplicates and semantically similar actions, we encountered issues like typographical errors, overly complex descriptions, and significant ambiguities in the remaining texts. These problems rendered the descriptions unsuitable for generating actionable human motion skills despite their uniqueness.

To further refine our text instructions, we evaluated the remaining descriptions for their suitability in model processing and practical motion generation. Our evaluation, detailed in Fig. A1, focused on fluency, conciseness, and the specificity of individual human poses within a brief sequence of frames. Descriptions that were direct and descriptive, containing clear verbs and nouns, were preferred over those with a sequential or ambiguous nature. Using a standardized scoring process, we ranked the action descriptions by their scores. After addressing issues in an initial round of scoring, a second evaluation was conducted to fine-tune our selection, as mentioned in Fig. A2. This led to the exclusion

### Score Prompt I

You are a language expert. Please rate the following actions on a scale of 0 to 10 based on their use of language. The requirements are:

1. *The description should be fluent and concise.*
2. *The description should correspond to a single human pose, instead of a range of possible poses.*
3. *The description should describe a human pose at a short sequence of frames instead of a long sequence of frames (this requirement is not mandatory).*
4. *If the description contains sequential logic, rate it lower. "Walk in a circle" is a kind of sequential logic.*
5. *Except for the subject, the description should have only one verb and one noun.*
6. *If the description is vivid (like "dances like Michael Jackson"), rate it higher.*

Here are some examples you graded in the last round:

- 6 - A person is swimming with his arms.
- 3 - Sway your hips from side to side.
- 7 - A person smashed a tennis ball.
- 4 - A person is in the process of sitting down.
- 5 - A person brings up both hands to eye level.
- 9 - A person dances like Michael Jackson.
- 2 - A person packs food in the fridge.
- 5 - A person flips both arms up and down.
- 8 - Looks like disco dancing.
- 3 - Kneeling person stands up.
- 1 - A person does a gesture while doing kudo.
- 6 - A person unzipping pants flyer.
- 0 - then kneels on both knees on the floor.
- 2 - A person is playing pitch and catch.
- 1 - A person gesturing them walking backward.
- 4 - A person seems confident and aggressive.
- 1 - A person circles around with both arms out.
- 5 - A person prepares to take a long jump.
- 6 - A person jumps twice into the air.
- 0 - Turning around and walking back.

Now, please provide your actions in the format 'x - yyyy,' where 'x' is the score, and 'yyyy' is the original sentence. Please note that Do not change the original sentence.

Figure A1. **Score Prompt I.** This prompt focuses on filtering text descriptions for fluency, conciseness, and specificity, particularly targeting individual human poses within a short sequence of frames.

of descriptions within certain score ranges (0-0.92, 0.98-0.99), resulting in a curated dataset of 1,896 unique action descriptions optimized for model training.

## Score Prompt II

You are a language expert. Please rate the following actions on a scale of 0 to 10 based on the ambiguity of the description. Examine whether this action description corresponds to a unique action. If the description corresponds to fewer actions, like "wave with both arms", rate it higher. If the description corresponds to abundant actions, like "do yoga", rate it lower.

- 7 - grab items with their left hand.
- 8 - hold onto a handrail.
- 9 - do star jumps.
- 5 - arms slightly curled go from right to left.
- 3 - sit down on something.
- 9 - kick with the right foot.
- 7 - stand and put arms up.
- 9 - cover the mouth with the hand.
- 8 - stand and salute someone.
- 2 - break dance.
- 6 - spin body very fast.
- 7 - open bottle and drink it.
- 2 - do the cha-cha.
- 5 - do sit-ups.
- 4 - slowly stretch.
- 6 - cross a high obstacle.
- 7 - grab something and shake it.
- 4 - lift weights to get buff.
- 8 - move left hand upward.
- 7 - walk forward swiftly.

Now, please provide your actions in the format 'x - yyyy,' where 'x' is the score, and 'yyyy' is the original sentence. Please note that Do not change the original sentence.

Figure A2. **Score Prompt II.** This prompt selects for direct and richly detailed action descriptions, prioritizing clarity with a distinct verb and noun over descriptions based on sequential or complex logic.

**Rewrite text** In the final refinement phase, we address the specificity of action descriptions, crucial for accurately generating motions. Vague descriptions, such as 'jump rope', can lead to ambiguous interpretations and various motion realizations, challenging the model's training due to the similarity of rewards for different motions. This observation is consistent with other motion generation studies utilizing CLIP [11, 43].

To enhance the clarity and effectiveness of the reward calculation, we rephrase and detail the descriptions. For instance, 'jump rope' is clarified to 'swinging a rope around your body', with further details like 'Raise both hands and shake them continuously while simultaneously jumping up

with both feet, repeating this cycle'. Additionally, we break down actions into more discrete moments, such as 'legs off the ground, wave hand', to improve the reward function's precision. Our methodology for this textual refinement is detailed in Fig. A3.

## Rewrite Prompt

Describe an action of instruction for a humanoid agent. The description must satisfy the following conditions:

1. The description should be concise.
2. The description should describe a human pose in a single frame instead of a sequence of frames.
3. The description should correspond to only one human pose, instead of a range of possible poses, minimize ambiguity.
4. The description should be less than 8 words.
5. The description should not contain a subject like "An agent", "A human".
6. The description should have less than two verbs and two nouns.
7. The description should not have any adjectives, adverbs, or any similar words like "with respect".
8. The description should not include details describing expressions or fingers and toes.

For example, it's better to describe "take a bow" as "bow at a right angle."

Figure A3. **Rewrite Prompt.** This prompt is designed for rephrasing action descriptions to enhance clarity and incorporate additional details, aiming to improve the specificity and effectiveness of the generated motions.

## A.2. Motion Data

For the study, we curated 93 motion clips, organizing them by movement type and style into a structured dataset. We delineated movements into three categories: *move\_around*, *act\_in\_place*, and *combined*; and styles into five categories: *attack*, *crawl*, *jump*, *dance*, and *usual*. The clips were then classified into these eight categories, with a weighting system applied based on the inverse frequency of each category to enhance the representation of less common actions. For motions that spanned multiple categories, their weights were averaged based on their inverse frequency values. This approach aimed to ensure a balanced action distribution within the dataset, emphasizing the inclusion of rarer actions to avoid overrepresentation of any single action type. The categorization and its impact on the dataset distribution are illustrated in the diagram available in Fig. A10.

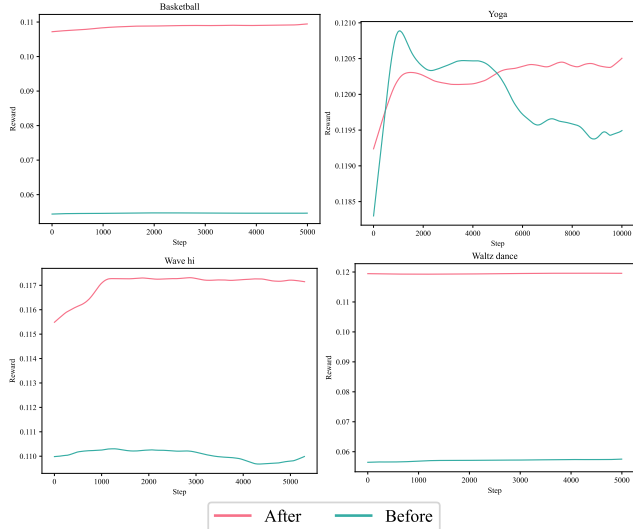


Figure A4. **Rewards before and after text enhancement.** The red curve depicts reward trends following text enhancement, contrasting with the pre-enhancement trends shown by the green curve.

## B. Experiments

This supplementary section expands on the experimental analyses from Sec. 4, focusing on the text description. Beyond the quantitative metrics addressed in the main document, we explore the changes in reward function dynamics pre- and post-text refinement across various instructions. This includes a detailed comparison of CLIP similarity scores during training to critically evaluate the effectiveness and design of different reward functions.

### B.1. Text Enhancement

Utilizing the text enhancement strategy described in Appendix A.1, we have refined action descriptions from existing open-source datasets, reducing ambiguity and enhancing clarity and applicability. To gauge the impact of these refined descriptions on training efficacy, we track and compare the reward feedback during the training phases.

Selecting four instructions at random from our dataset for illustration, we compare reward trends before and after text enhancements—represented by green and red curves, respectively, in our graphs. This comparison reveals that refined instructions consistently yield superior reward trajectories from the start, showing a swift and steady ascent to a performance plateau. This indicates that text enhancement notably improves policy training efficiency and convergence speed. Specifically, for intricate actions like *Yoga* (as shown in the top right figure of Fig. A4), refined instructions result in a more stable and gradual reward increase, signifying improved training stability and model performance.

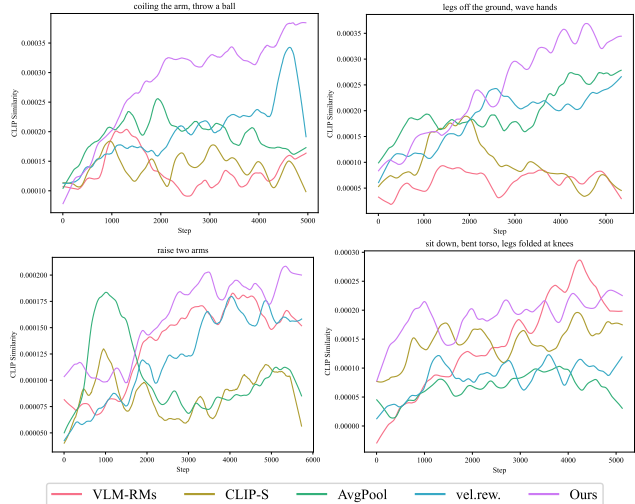


Figure A5. **The CLIP similarity calculated by different reward designs.**

### B.2. Implementation Details

### B.3. Reward Function Analysis

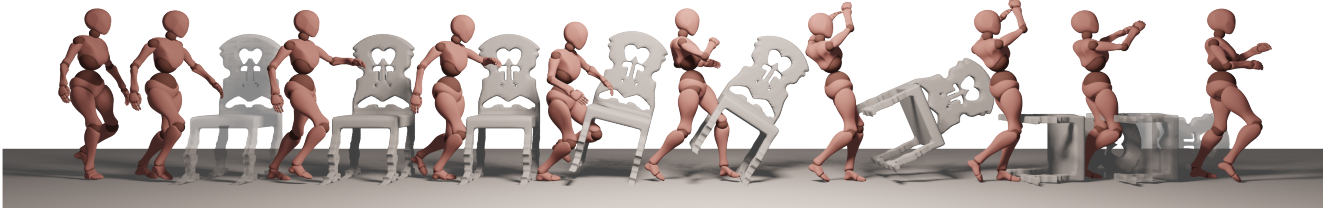
To evaluate and compare various reward function designs, we use cosine similarity between image and text features as a uniform metric, accommodating the differing numerical scales inherent to each reward design. As depicted in Fig. A5, we represent five reward functions using distinct colors, with our method marked in purple.

Aligning with discussions in the main text (Fig. 5), we examine four instructions from our user study for a detailed comparison. Our findings indicate that our method uniformly improves image-text alignment throughout training, achieving consistent convergence. While some methods exhibit comparable performance on select instructions, they generally show less consistency, with initial gains often receding over time. In contrast, our approach demonstrates robustness against the variabilities of open-vocabulary training, leading to stable and reliable performance improvements.

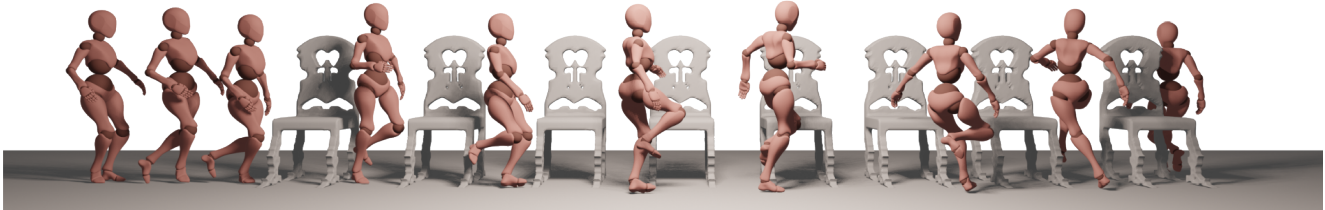
To assist readers in replicating our work, we have included a comprehensive breakdown of hyperparameter settings in Tabs. A1 and A2.

### B.4. Interaction Motions

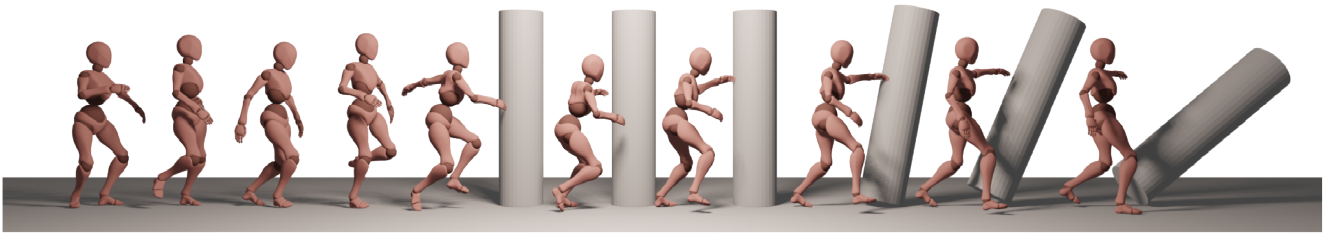
Within the main text, we highlighted AnySkill’s proficiency in mastering tasks involving interactions with diverse objects, underscoring its capability to adapt across a spectrum of interaction scenarios. For experimental validation, we deliberately chose a range of objects, both rigid (*e.g.*, pillars, balls) and articulated (*e.g.*, doors, chairs), to demonstrate the method’s versatility. The quantitative analyses of these object interactions, as detailed in Appendix B.2, affirm the flexibility of our approach. Our system is shown to adeptly navigate a variety of action requirements, as speci-



(a) kick the white chair



(b) move around the white chair



(c) strike the pillar

Figure A6. Additional results of interaction motions.

Table A1. Hyperparameters used for the training of low-level controller.

Hyper-Parameters	Values
dim(Z) Latent Space Dimension	64
Encoder Align Loss Weight	1
Encoder Uniform Loss Weight	0.5
$w_{gp}$ Gradient Penalty Weight	5
Encoder Regularization Coefficient	0.1
Samples Per Update Iteration	131072
Policy/Value Function Minibatch Size	16384
Discriminators/Encoder Minibatch Size	4096
$\gamma$ Discount	0.99
Learning Rate	$2 \times 10^{-5}$
GAE( $\lambda$ )	0.95
TD( $\lambda$ )	0.95
PPO Clip Threshold	0.2
$T$ Episode Length	300

fied by different text descriptions, maintaining efficacy even when faced with repetitive initial conditions or identical objects.

Table A2. Hyperparameters used for the training of high-level controller.

Hyper-Parameters	Values
$w_{gp}$ Gradient Penalty Weight	5
Encoder Regularization Coefficient	0.1
Samples Per Update Iteration	131072
Policy/Value Function Minibatch Size	16384
Discriminators/Encoder Minibatch Size	4096
$\gamma$ Discount	0.99
Learning Rate	$2 \times 10^{-5}$
GAE( $\lambda$ )	0.95
TD( $\lambda$ )	0.95
PPO Clip Threshold	0.2
$T$ Episode Length	300

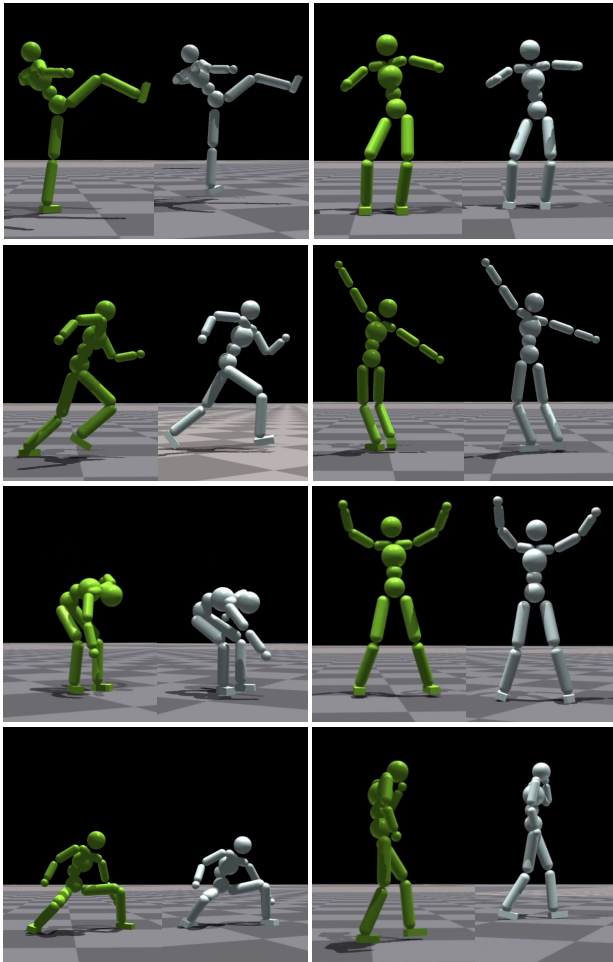


Figure A7. **Atomic actions from the trained low-level controller.** In each subfigure, the green agent shows the reference motion from the dataset, and the white agent shows our learned atomic action.



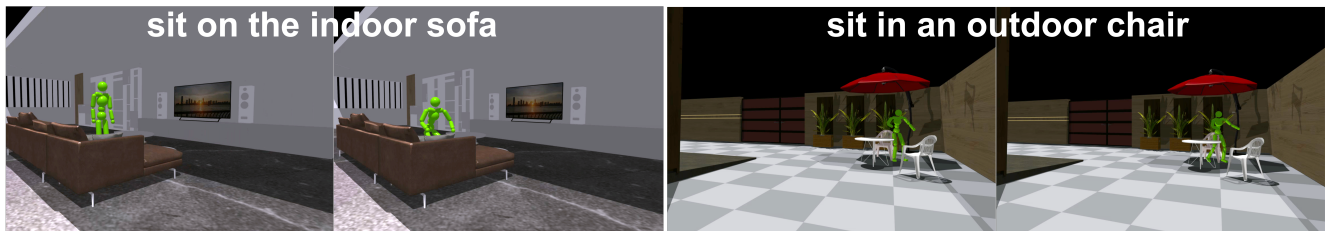
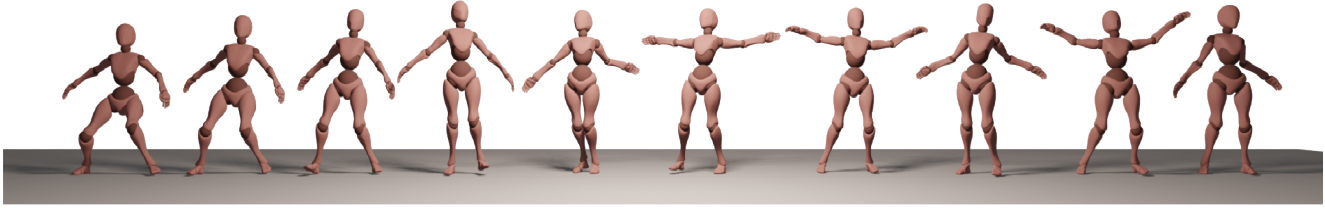
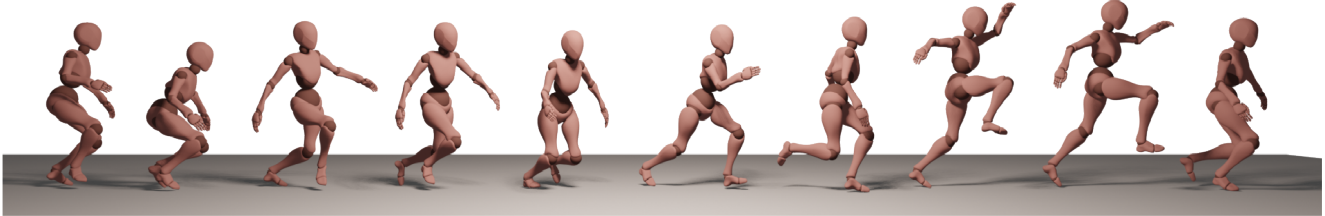


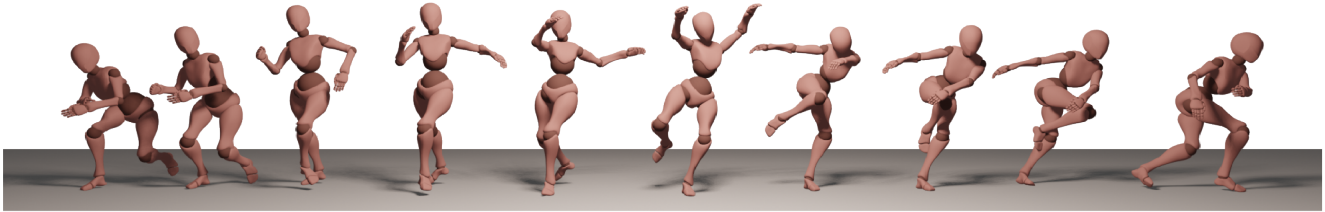
Figure A8. **Real-time scene interaction.** We employed both indoor and outdoor scenes within IsaacGYM. Throughout the training process, we conducted real-time rendering and obtained feedback on physical interactions.



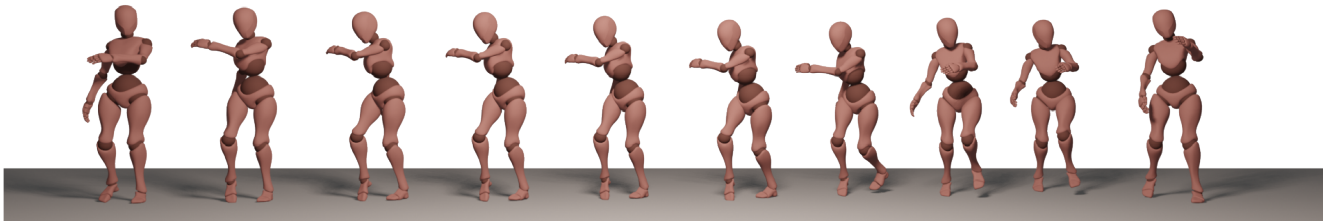
(a) wave hands up and down



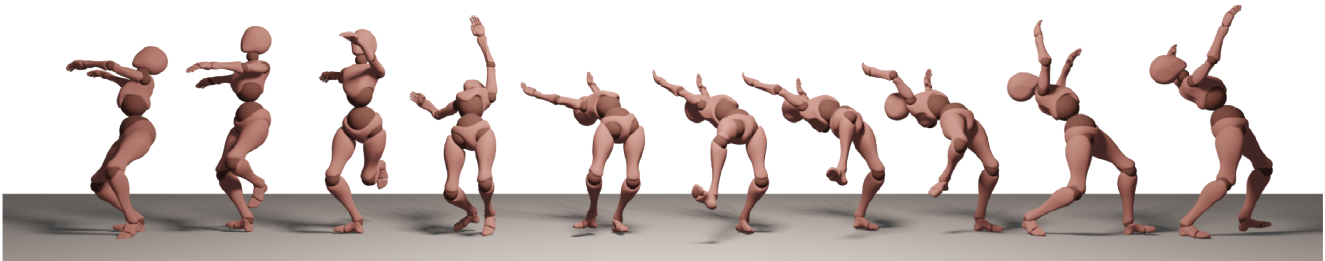
(b) jump high



(c) left leg forward, right leg retreats



(d) raise one arm, put the other hand down



(e) raise hands above head, bend body



(f) hit a tennis smash with arm

Figure A9. More results of open-vocabulary physical skills.

