



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

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

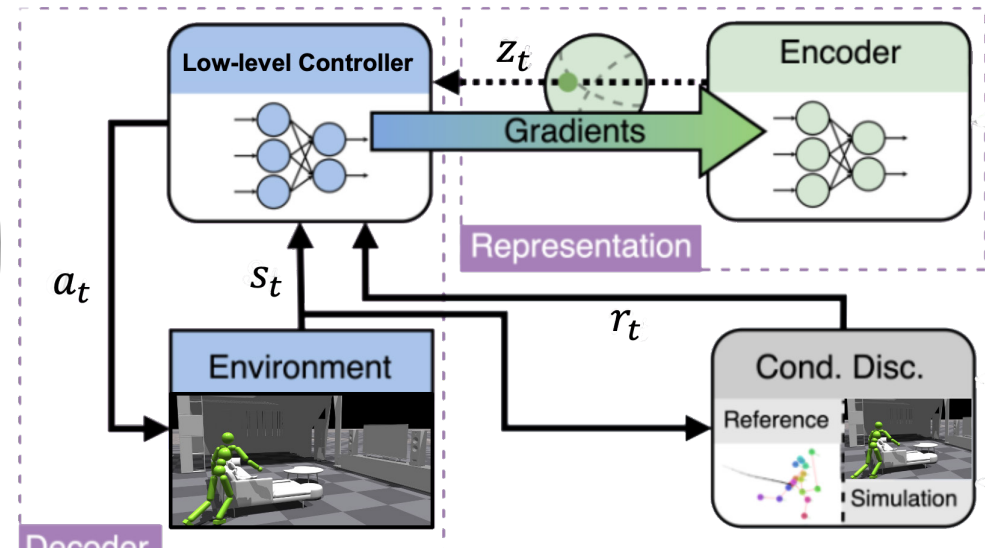
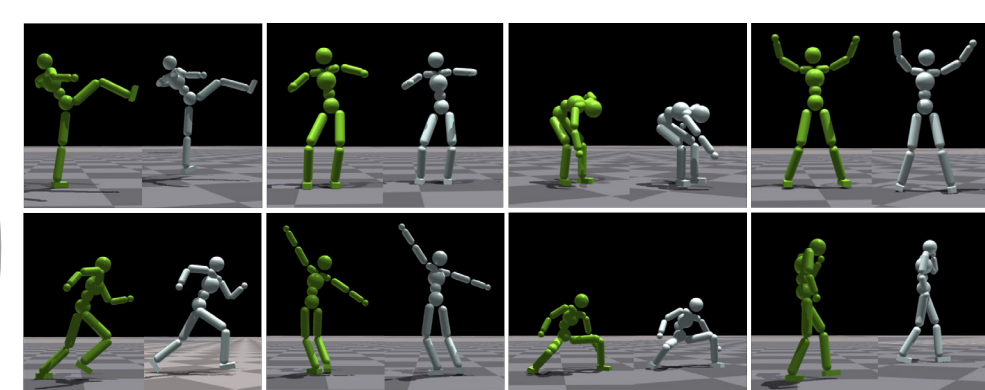
<sup>1</sup>Institute for Artificial Intelligence, Peking University, <sup>2</sup>State Key Laboratory of General Artificial Intelligence, BIGAI



<https://anyskill.github.io/>

## Low-level Algorithm

- ✓ **93** modified reference motions
- ✓ **Semantic** latent space



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

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

## High-level Algorithm

- ✓ **Only** image-based reward
- ✓ **Agent-centered** image rendering
- ✓ **LLM** enhances open-vocabulary capabilities.

```

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  $z$  from  $\mathcal{Z}$ 
6     if horizon_length = 1 then
7        $s \leftarrow$  initialize;  $z \leftarrow z$ 
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         $s, p \leftarrow$  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
  
```

## ANY Text Description

- a. Sit down
- b. Dance
- c. Kick
- d. Jump rope

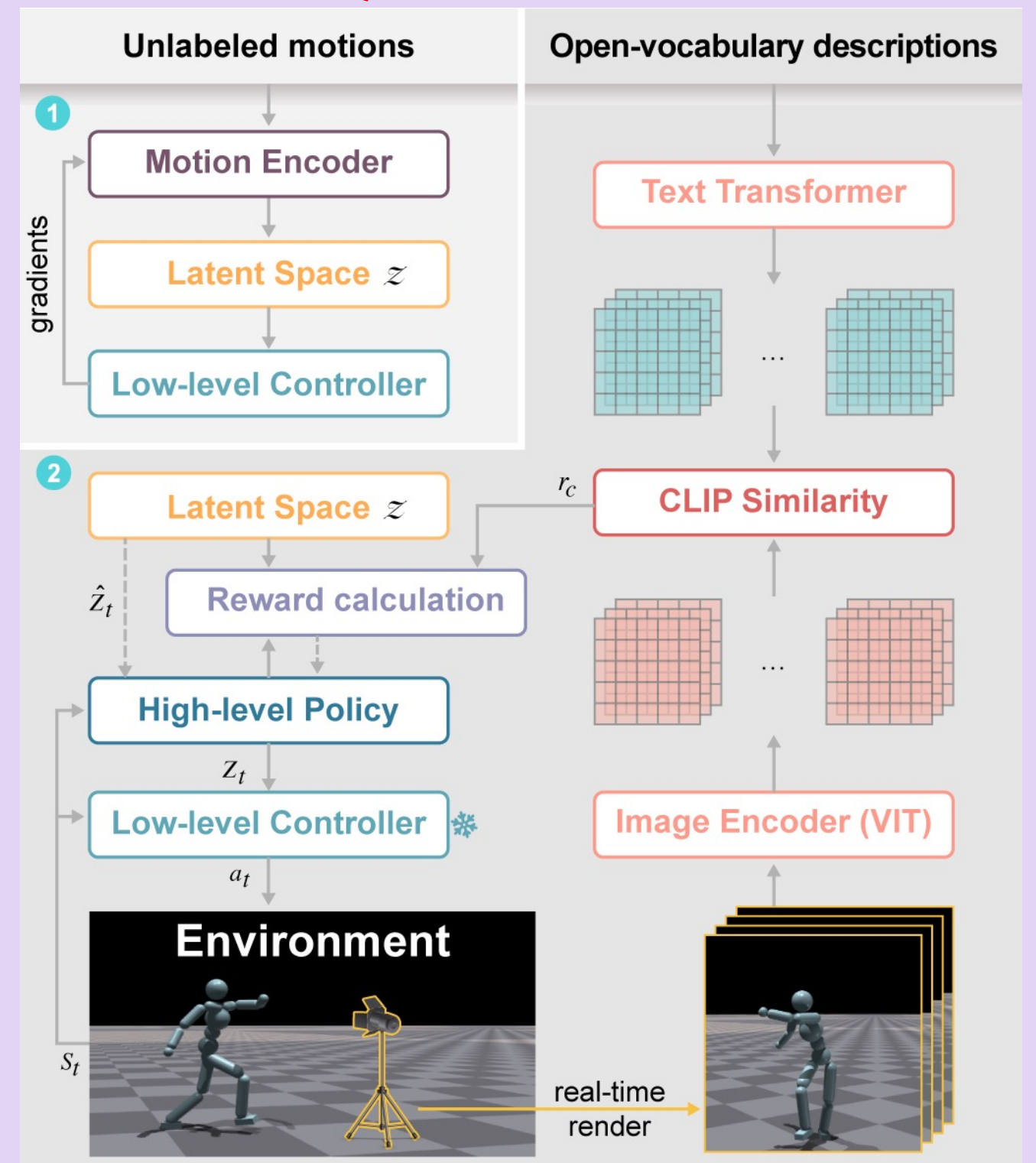


### Iterate text enhancement

1. Describe one specific human pose.
2. Ensure the verb and noun relate to a human limb.
3. Make the description fluent, concise, and unambiguous. ...

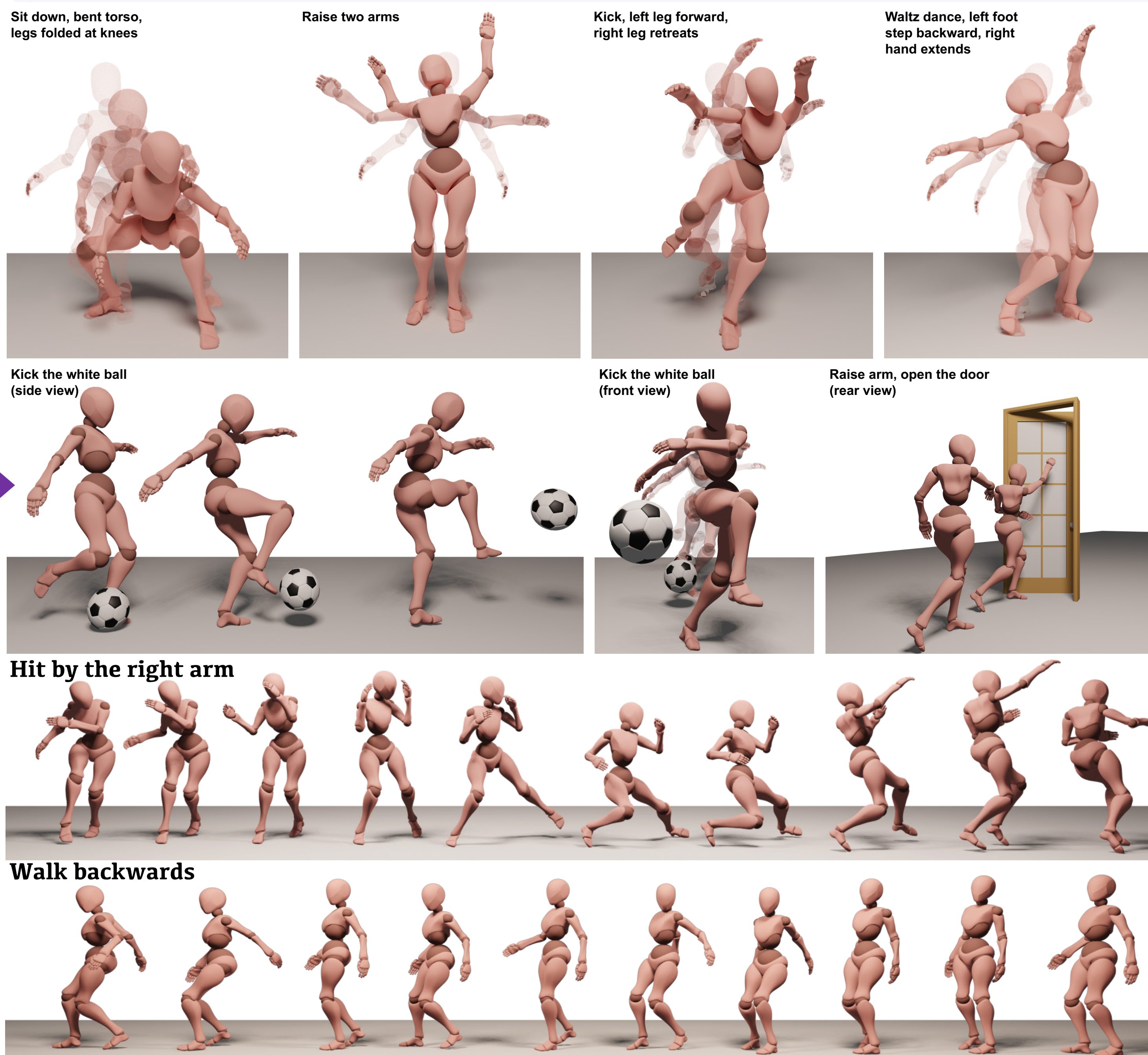
- a. Sit down, bent torso, legs folded at knees
- b. Dance, left foot step backward, right hand extends
- c. Kick, left leg forward, right leg retreats
- d. Jump rope, legs off the ground, wave hands

## Model w/o manual reward



## Contributions

- ✓ AnySkill is a **hierarchical approach** designed for learning open-vocabulary physical skills.
- ✓ We create flexible, generalizable, **image-based** rewards, **eliminating** the need for **manual** design.
- ✓ Our method outperforms existing approaches and enables agents to **interact smoothly with dynamic objects** in various contexts.



## Interaction with object & scene



## Quantitative Results

Table 1. Quantitative evaluation of high-level policy.

	Success $\uparrow$	Natural $\uparrow$	Smooth $\uparrow$	Physics $\uparrow$	CLIP-S $\uparrow$
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 [56])	2.37	4.90	5.65	6.41	21.35
Ours (full)	<b>6.16</b>	<b>6.23</b>	<b>6.51</b>	<b>6.93</b>	<b>24.18</b>

Table 2. Quantitative evaluation of interaction motions.

	Success $\uparrow$	Natural $\uparrow$	Smooth $\uparrow$	Physics $\uparrow$	CLIP-S $\uparrow$
Interaction w. object	5.42	-0.74	5.62	-0.61	5.34
Interaction w. scene	4.53	-1.63	4.47	-1.76	5.01
Ours	5.41	-1.52	22.41	-1.73	

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 [69]	3.80	5.41	5.98	6.21	19.78
AvgPool [56]	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
Ours	<b>6.16</b>	<b>6.23</b>	6.51	<b>6.93</b>	<b>24.18</b>