

## Appendix

### A. Reward Function

Reward items are listed in TABLE III. The task reward helps the robot adapt to the specific task scenarios, and the details of the design (4,6,7) is shown in II.B above. Where we set the Goal Velocity’s reward coefficient  $k_{positive}$  to 1 and the penalty coefficient  $k_{negative}$  to -3. This set of parameters was optimized based on our comprehensive simulation experiments. Additionally, we designed the regularization reward referring to [46], adjusted it according to the type of robots and task scenarios.

TABLE III: Reward Function

Term	Equation	Weight
Task reward		
Goal Velocity	$r_{vel}(4)$	1.5
Yaw Angular Velocity	$exp(-4 \omega_{yaw}^{cmd} - \omega_{yaw} )$	0.5
Hip Position	$r_{Pos}(6)$	-0.5( $\omega_1=\omega_2=0.5$ )
Collision	$r_{collision}(7)$	-10.0
Regularization reward		
Z Velocity	$v_z^2$	-0.5
X&Y Velocity	$\omega_x^2 + \omega_y^2$	-0.01
Dof Acceleration	$\sum_{i=1}^{12} \ddot{q}_i^2$	$-2.5 * 10^{-7}$
Action Rate	$\sqrt{\sum_{i=1}^{12} (a_t - a_{t-1})^2}$	-0.1
Delta Torques	$\sum_{i=1}^{12} (\tau_t - \tau_{t-1})^2$	$-1.0 * 10^{-7}$
Torques	$\sum_{i=1}^{12} (\tau_t)^2$	$-1.0 * 10^{-5}$
Dof Error	$\sum_{i=1}^{12} (q - q_{default})^2$	-0.04
Feet Stumble	$ F_{feet}^{hor}  > 4 *  F_{feet}^{ver} $	-1
Dof Position Limits	$\sum_{i=1}^{12} (q_i^{out} \text{ if } q_i > q_{max} \text{ or } q_i < q_{min})$	-10.0

### B. Additional Training Details

**Network Architecture** Our learning framework’s overall network architecture includes not only the components shown in Figure 2 but also some additional modules not depicted in the figure. The teacher policy consists of six Multilayer Perceptron (MLP) parts: collision domain encoder, collision estimator, privileged information encoder, privileged information estimator, velocity estimator, and teacher policy network. Where  $h_t$  denotes collision domain information and  $g_t$  denotes privilege information. The privileged information encoder supervises the privileged information estimator and optimizes learning using the ROA method. Privileged Estimator’s network type is also a CNN with a similar structure to the Collision Estimator. The student policy includes five parts: the collision estimator, velocity estimator, and privileged information estimator from the teacher policy, along with the hybrid imagination model comprised of a Gated Recurrent Unit (GRU) network and an MLP-based student policy network. Table I provides more details on each layer.

**Training course** Course learning is crucial for robots to effectively travel obstacles in complex environments. Without this capability, robots would struggle to learn effectively. We have implemented the velocity-travel course training

TABLE IV: Network architectures

Module	Inputs	Hidden Layers	Outputs
Teacher policy			
Collision Domain Enc	$h_t$	[256, 128, 64]	$p_t$
Collision Estimator	$o_{t-10}, \dots, o_{t-1}, o_t$	/	$\hat{c}_t$
Privileged Encoder	$g_t$	[64, 32]	$e_t$
Privileged Estimator	$o_{t-5}, \dots, o_{t-1}, o_t$	/	$\hat{e}_t$
Velocity Estimator	$o_t$	[128, 64]	$\hat{v}_t$
Teacher Network	$o_t, p_t, \hat{c}_t, e_t, \hat{v}_t$	[512, 256, 128]	$a_t$
Student policy			
Hybrid Imagination Model	$o_{t-10}, \hat{c}_{t-10}, \dots, \hat{c}_t, o_t$	/	$\hat{p}_t$
Student Network	$o_t, \hat{p}_t, \hat{c}_t, \hat{e}_t, \hat{v}_t$	[512, 256, 128]	$\hat{a}_t$

method[43], [46], where the robot’s linear velocity is randomly sampled within the range of [0, 1]. If the robot’s travel distance in one iteration exceeds half of the preset heading speed integral, the terrain difficulty is increased; otherwise, it is decreased. The levels of terrain difficulty are detailed in TABLE I.

Hyperparameter	Value
Discount Factor	0.99
GAE Parameter	0.95
Timesteps per Rollout	21
Epochs per Rollout	5
Minibatches per Epoch	4
Entropy Bonus ( $\alpha_2$ )	0.01
Value Loss Coefficient ( $\alpha_1$ )	1.0
Clip Range	0.2
Reward Normalization	yes
Learning Rate	2e-4
# Environments	4096
Optimizer	Adam

TABLE V: PPO hyperparameters.

**Policy training and Imitation training** We use PPO with hyperparameters listed in TABLE V to train the teacher policy. We regard the process of teacher policy supervising student policy learning as imitation learning process,  $a_t$  and  $\hat{a}_t$  are the action vectors from the teacher and student respectively in an actor network. The gradient of the loss  $\mathcal{L}$  can be defined as the sum of the squared norms of the difference between  $a_t$  and  $\hat{a}_t$  over all actions:

$$\mathcal{L}_{Imitation} = \|\pi_{teacher}(\cdot|s) - \pi_{student}(\cdot|s)\|_2^2$$

Here,  $\|\cdot\|_2$  denotes the  $L^2$  norm. Where  $\pi_{teacher}(\cdot|s)$  is the action from the Teacher policy and  $\pi_{student}(\cdot|s)$  is the action from the Student policy .