

# 1 Appendix

## 2 1.1 Detailed Implementation of Real Robot System

3 **Bimanual Hands System.** For our system, we have developed a ROS-based pipeline that operates at  
4 a control frequency of 20Hz. This pipeline serves as the foundation for controlling our setup, enabling  
5 efficient communication and coordination between the different components. In our configuration, the  
6 Arm-Hand subsystems are controlled by a single policy utilizing multiple agents. This unified policy  
7 governs the actions of both subsystems, promoting synchronized and collaborative behavior in our  
8 setup. To achieve this, we control the motion of the robotic arms through Modbus TCP (Transmission  
9 Control Protocol) using an AC/DC Control Box. The control boxes of the two robotic arms are  
10 connected to a router via Ethernet cables, and the router is then connected to the host computer.  
11 Additionally, the two robot hands are directly connected to the same computer using RS-485 serial  
12 communication.

13 **Object Tracking.** Real-time object tracking is performed with an Intel RealSense D435 stereo  
14 camera. Since the object has a high color contrast from its background, we first use a simple color  
15 detector on the RGB image to find the pixel location of the object. The color range for detecting a  
16 blue object is constrained between  $[80, 200, 0]$  and  $[120, 255, 0]$  in HSV color space. Next, the 3D  
17 position of that pixel is obtained from querying the corresponding depth value on the depth image,  
18 where post-processing filters including disparity, spatial and temporal, are applied to reduce depth  
19 noise. Finally, we get the 3D object position in robot frame from image frame with calibrated camera  
20 extrinsics parameters.

## 21 1.2 Sim2Real Transfer

22 **System Identification.** To achieve a successful sim-to-real transfer, we utilize system identification  
23 techniques to align the behavior of the PD (Proportional-Derivative) controller of the arm and hand  
24 in simulation with that in the real world. This involves tuning the PD coefficients of the controllers  
25 to ensure that their responses to impulse and sinusoidal inputs are aligned. This step is crucial in  
26 ensuring that the control actions generated in simulation can be effectively applied to the real robot  
27 setup, enabling a reliable sim-to-real transfer of our system.

## 28 Domain Randomization

29 Isaac Gym offers several domain randomization functions for reinforcement learning training. We  
30 apply randomization to the task, as indicated in Table. 1 for each environment. We generate new  
31 randomizations every 1000 simulation steps.

Table 1: Domain randomization parameters.

Parameter	Type	Distribution	Initial Range
<b>Robot</b>			
Mass	Scaling	uniform	[0.5, 1.5]
Friction	Scaling	uniform	[0.7, 1.3]
Joint Lower Limit	Scaling	loguniform	[0.0, 0.01]
Joint Upper Limit	Scaling	loguniform	[0.0, 0.01]
Joint Stiffness	Scaling	loguniform	[0.0, 0.01]
Joint Damping	Scaling	loguniform	[0.0, 0.01]
<b>Object</b>			
Mass	Scaling	uniform	[0.5, 1.5]
Friction	Scaling	uniform	[0.5, 1.5]
Scale	Scaling	uniform	[0.95, 1.05]
<b>Observation</b>			
Obs Correlated. Noise	Additive	gaussian	[0.0, 0.001]
Obs Uncorrelated. Noise	Additive	gaussian	[0.0, 0.002]
<b>Action</b>			
Action Correlated Noise	Additive	gaussian	[0.0, 0.015]
Action Uncorrelated Noise	Additive	gaussian	[0.0, 0.05]
<b>Environment</b>			
Gravity	Additive	normal	[0, 0.4]

32 **1.3 Hyperparameters of the RL algorithms**

Table 2: Hyperparameters of MAPPO.

Hyperparameters	Throw and Catch
Num mini-batches	1
Num opt-epochs	5
Num episode-length	8
Hidden size	[1024, 1024, 512]
Use popart	True
Use value norm	True
Use proper time limits	False
Use huber loss	True
Huber delta	10
Clip range	0.2
Max grad norm	10
Learning rate	5.e-4
Opt-eps	5.e-4
Discount ( $\gamma$ )	0.96
GAE lambda ( $\lambda$ )	0.95
Std x coef	1
Std y coef	0.5
Ent-coef	0

Table 3: Hyperparameters of PPO.

Hyperparameters	Throw and Catch
Num mini-batches	1
Num opt-epochs	5
Num episode-length	8
Hidden size	[1024, 1024, 512]
Use popart	True
Use value norm	True
Use proper time limits	False
Use huber loss	True
Huber delta	10
Clip range	0.2
Max grad norm	10
Learning rate	5.e-4
Opt-eps	5.e-4
Discount ( $\gamma$ )	0.96
GAE lambda ( $\lambda$ )	0.95
Std x coef	1
Std y coef	0.5
Ent-coef	0