

# HT\_Pi Reinforcement Learning Baseline

## Technical Document

This document introduces a reinforcement learning environment based on NVIDIA Isaac Gym. Pi\_rl\_baseline is a sim2sim framework for High Torque Electromechanical's bipedal robot Mini Pi, supporting policy transfer from Isaac Gym to Mujoco. This framework enables users to verify the trained policies in different physics engines, thereby evaluating their robustness and generalization ability.

It is recommended to learn alongside the README of High Torque's open-source project ([https://github.com/HighTorque-Robotics/livelybot\\_pi\\_rl\\_baseline](https://github.com/HighTorque-Robotics/livelybot_pi_rl_baseline)) for a more comprehensive understanding.

The following is the introductory tutorial for reinforcement learning by High Torque Electromechanical, suitable for demonstrating to beginners how to train and transfer policies based on a simulation environment:

1. Install Miniconda or Anaconda independently;
2. Create a virtual environment:

代码块

```
1 conda create -n pi_env python=3.8
```

3. Check the NVIDIA graphics card driver:

Use the command `nvidia-smi` in the command line to check the CUDA version of the driver. You can see the CUDA version is 12.8 and the driver version is 570.

NVIDIA-SMI 570.133.20			Driver Version: 570.133.20		CUDA Version: 12.8	
GPU	Name	Persistence-M	Bus-Id	Disp.A	Volatile	Uncorr. ECC
Fan	Temp	Perf	Pwr:Usage/Cap	Memory-Usage	GPU-Util	Compute M.
						MIG M.
0	NVIDIA GeForce RTX 4060 Ti	Off	00000000:01:00.0	On		N/A
0%	42C	P8	12W / 165W	1067MiB / 16380MiB	23%	Default
						N/A
Processes:						
GPU	GI	CI	PID	Type	Process name	GPU Memory Usage
	ID	ID				
0	N/A	N/A	1380	G	/usr/lib/xorg/Xorg	59MiB
0	N/A	N/A	2194	G	/usr/lib/xorg/Xorg	522MiB
0	N/A	N/A	2330	G	/usr/bin/gnome-shell	86MiB
0	N/A	N/A	2563	G	.../sunloginclient --cmd=autorun	12MiB
0	N/A	N/A	2707	G	...me/58.0.3029.81 Safari/537.36	2MiB
0	N/A	N/A	2733	G	...D4C2B51A7D651DD79862114E7657D	6MiB
0	N/A	N/A	3852	G	/usr/lib/firefox/firefox	193MiB
0	N/A	N/A	346176	G	...OTP --variations-seed-version	31MiB

```
conda install pi_env
```

```
sunteng@zixiang:~$ conda activate pi_env
(pi_env) sunteng@zixiang:~$
```

## 5. Install PyTorch:

## 代码块

```
1 conda install pytorch torchvision torchaudio pytorch-cuda=12.4 -c pytorch -c
nvidia
```

(The selected CUDA version must be less than or equal to the version installed on the computer; CUDA 12.4 is selected here.)

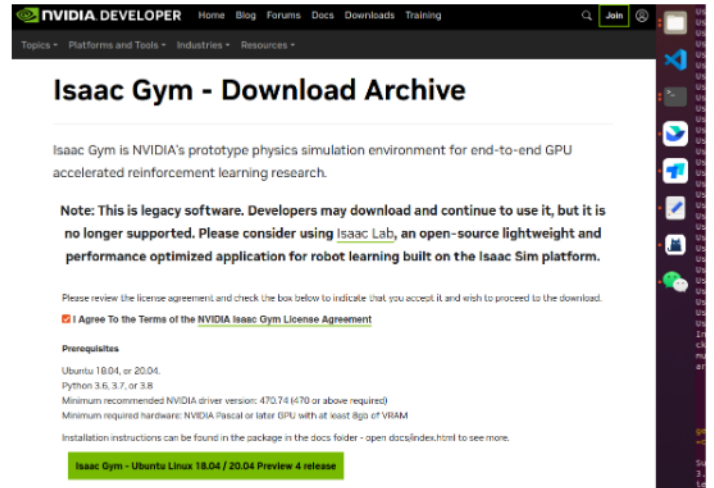
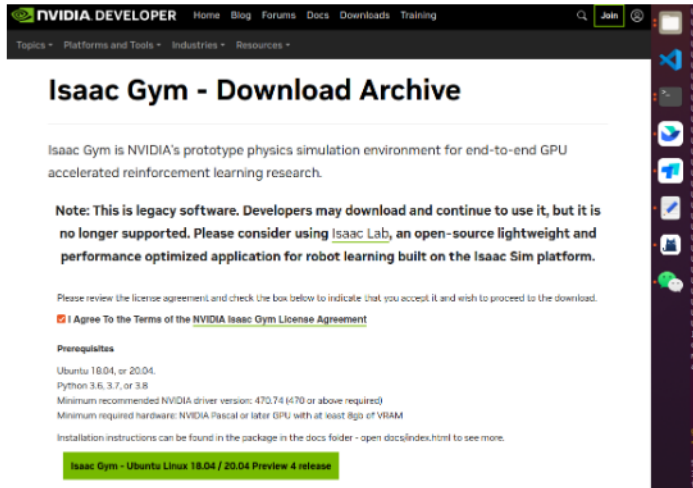
## 6. Install numpy using conda:

## 代码块

```
1 conda install numpy=1.23
```

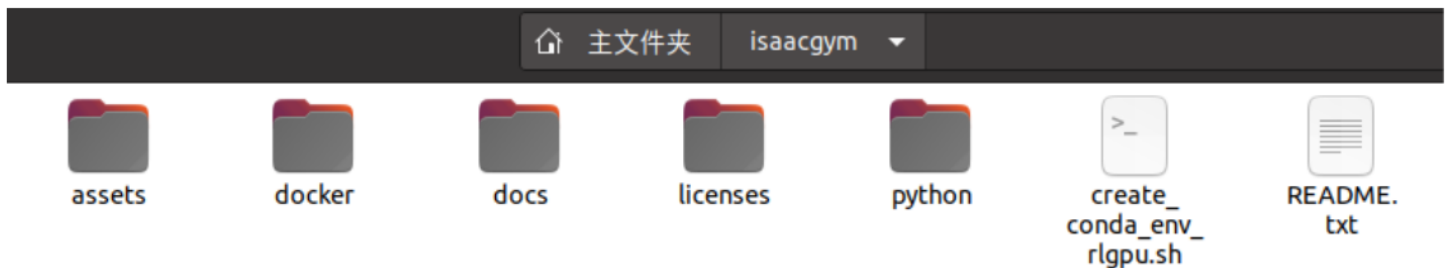
## 7. Install Isaac Gym:

Download and install Isaac Gym Preview 4 from the NVIDIA official website (Note: Only supported on Ubuntu 18.04 or 20.04).



## 8. Extract to the home directory:

Extract the downloaded Isaac\_Gym\_Preview\_4\_package to the home directory, and cut the isaacgym package under the subfolder to the home directory.



## 9. Enter the isaacgym package for installation:

First enter:

代码块

```
1 cd isaacgym/python
```

Then enter:

代码块

```
1 pip install -e .
```

```
(pi_env) sunteng@zixiang:~$ cd isaacgym/
(pi_env) sunteng@zixiang:~/isaacgym$ cd python/
(pi_env) sunteng@zixiang:~/isaacgym/python$ pip install -e .
```

## 10. Test if the environment installation is successful:

First enter:

代码块

```
1 cd examples
```

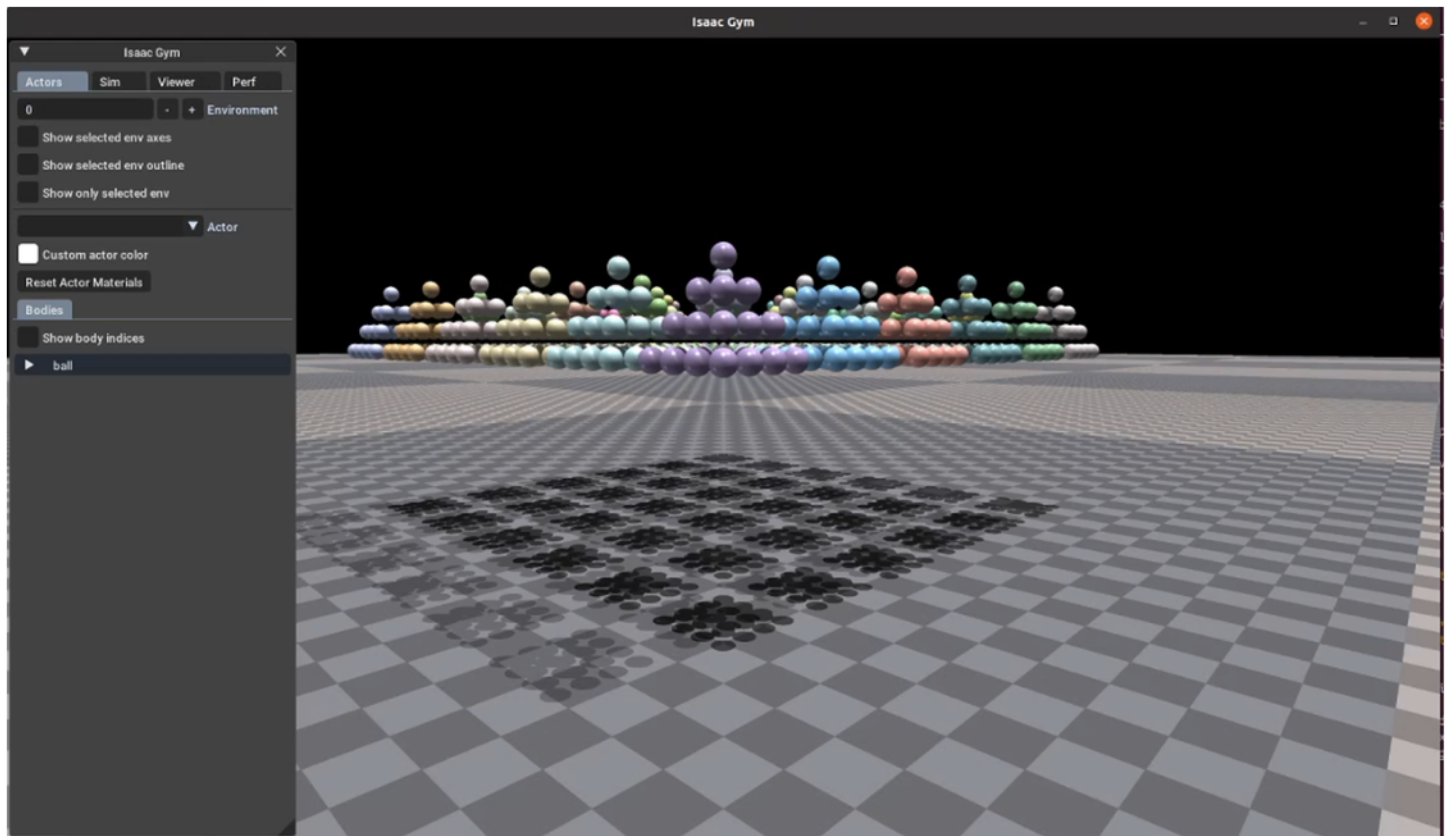
Then enter:

代码块

```
1 python 1080_balls_of_solitude.py
```

```
(pi_env) sunteng@zixiang:~/isaacgym/python$ cd examples/  
(pi_env) sunteng@zixiang:~/isaacgym/python/examples$ python 1080_balls_of_solitude.py
```

If the page shown in the figure below appears, the environment installation is successful.



## 11. Download the zip package of pi\_rl\_baseline:

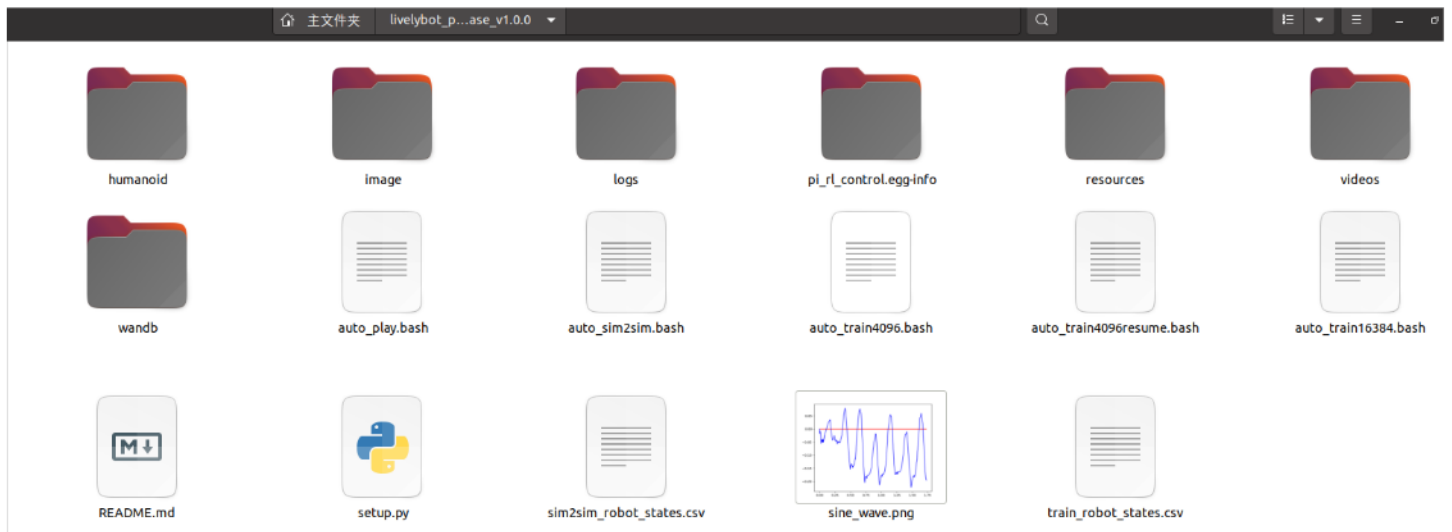
After successful download, extract it to the home directory. Enter

```
livelybot_pi_rl_baseline-release_v1.0.0(1)
```

 and cut

```
livelybot_pi_rl_baseline-release_v1.0.0
```

 to the home directory.



## 12. Install the current project in development mode:

- First enter:

代码块

```
1 cd livelybot_pi_rl_baseline-release_v1.0.0
```

- Then enter:

代码块

```
1 pip install -e .
```

(Note: The '.' must be included here.)

```
(pi_env) sunteng@zixiang:~$ cd livelybot_pi_rl_baseline-release_v1.0.0/
(pi_env) sunteng@zixiang:~/livelybot_pi_rl_baseline-release_v1.0.0$ pip install -e .
```

Finally, a prompt of "successfully installed" for a series of dependencies will appear, indicating successful configuration.

## 13. Train the PPO policy with 4096 environments and "v1" as the training version:

- First enter:

代码块

```
1 cd humanoid
```

- Then enter:

#### 代码块

```
1 python scripts/train.py --task=pai_ppo --run_name v1 --headless --num_envs 4096
```

If the interface shown in the figure below appears, it means the PPO policy training with "v1" as the training version is in progress.

```
wandb: Currently logged in as: 2742974891 (2742974891-beijing-union-university) to https://api.wandb.ai. Use 'wandb login --relogin' to force relogin
wandb: Tracking run with wandb version 0.19.11
wandb: Run data is saved locally in /home/sunteng/livelybot_pi_rl_baseline-release_v1.0.0/humanoid/wandb/run-20250603_150509-vpsmk5tb
wandb: Run 'wandb offline' to turn off syncing.
wandb: Syncing run Jun03_15-05-08_Pai_ppo_v1
wandb: ⭐ View project at https://wandb.ai/2742974891-beijing-union-university/pai_ppo
wandb: 🔍 View run at https://wandb.ai/2742974891-beijing-union-university/pai_ppo/runs/vpsmk5tb
wandb: ⚠️ WARNING Found log directory outside of given root_logdir, dropping given root_logdir for event file in /home/sunteng/livelybot_pi_rl_baseline-release_v1.0.0/logs/Pai_ppo/Jun03_15-05-08_v1
#####
Learning iteration 0/10001
Computation: 31499 steps/s (collection: 2.930s, learning 0.190s)
Value function loss: 0.3565
Surrogate loss: -0.0018
Mean action noise std: 1.00
Mean reward: 2.02
Mean episode length: 23.09
Mean episode rew_action_smoothness: -0.0009
Mean episode rew_base_acc: 0.0000
Mean episode rew_base_height: 0.0069
Mean episode rew_collision: 0.0000
Mean episode rew_default_ankle_roll_pos: 0.0044
Mean episode rew_default_hip_roll_joint_pos: 0.0195
Mean episode rew_default_thigh_joint_pos: 0.0064
Mean episode rew_dof_acc: -0.0002
Mean episode rew_dof_vel: -0.0002
Mean episode rew_feet_air_time: 0.0002
Mean episode rew_feet_clearance: 0.0139
Mean episode rew_feet_contact_forces: -0.0001
Mean episode rew_feet_contact_number: 0.0095
Mean episode rew_feet_distance: 0.0030
Mean episode rew_foot_slip: -0.0018
Mean episode rew_joint_pos: 0.0167
Mean episode rew_knee_distance: 0.0034
Mean episode rew_low_speed: -0.0007
Mean episode rew_orientation: 0.0021
Mean episode rew_termination: -0.0005
Mean episode rew_torques: -0.0000
Mean episode rew_track_vel_hard: -0.0029
Mean episode rew_tracking_ang_vel: 0.0051
Mean episode rew_tracking_lin_vel: 0.0180
Mean episode rew_vel_mismatch_exp: 0.0040
-----
Total timesteps: 98304
Iteration time: 3.12s
Total time: 3.12s
ETA: 31210.8s
```

Note: A wandb registration website will pop up for account registration (wandb is used in Isaac Gym for experiment tracking, hyperparameter recording, and visualization of reward curves for training metrics).

## 14. Evaluate the trained policy and automatically export a JIT model suitable for deployment:

- First install the onnx package:

#### 代码块

```
1 pip install onnx
```

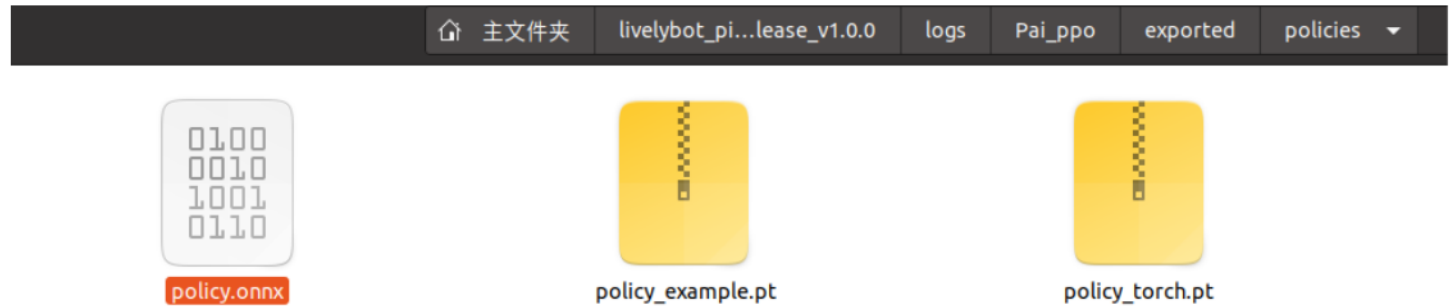
- Then enter:

#### 代码块

```
1 python scripts/play.py --task=pai_ppo --run_name v1
```



You can see `policy.onnx` under `/home/sunteng/livelybot_pi_rl_baseline-release_v1.0.0/logs/Pai_ppo/exported/policies`.

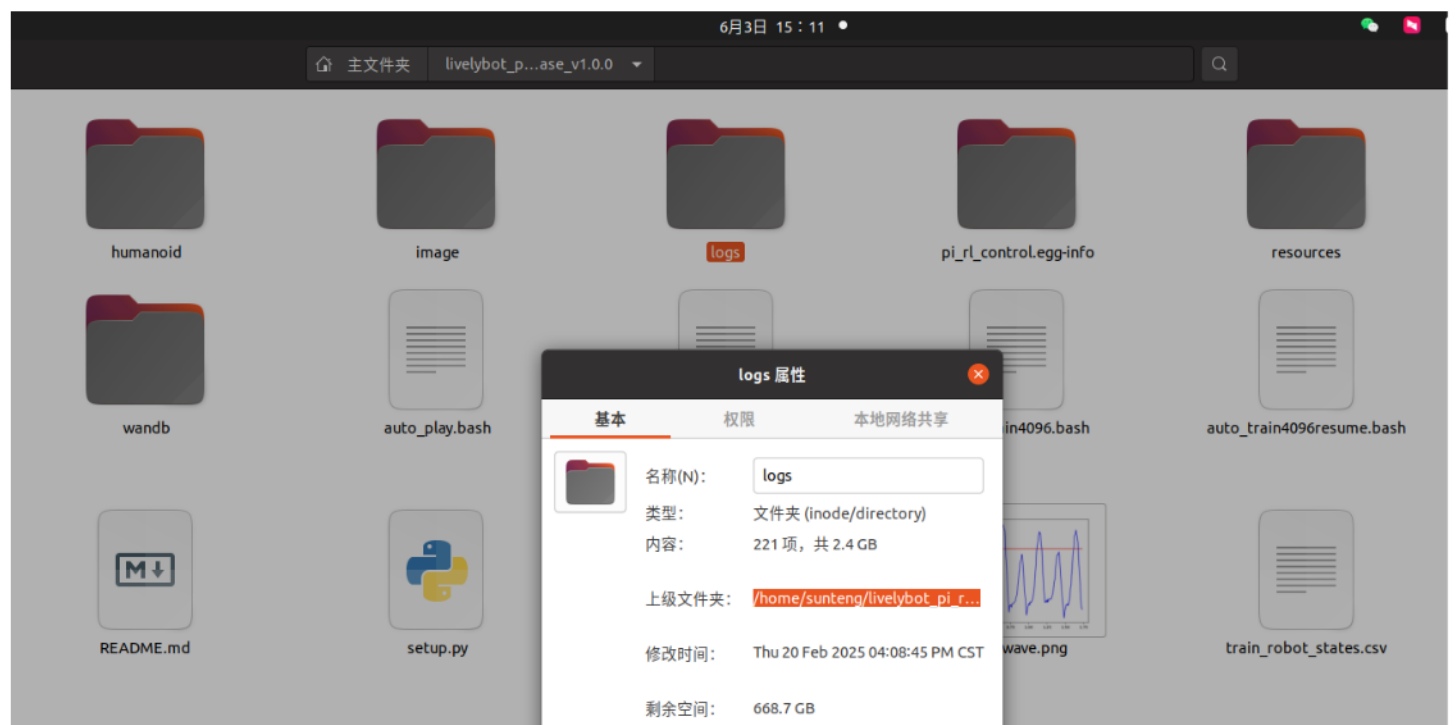


## 15. Implement sim2sim using Mujoco:

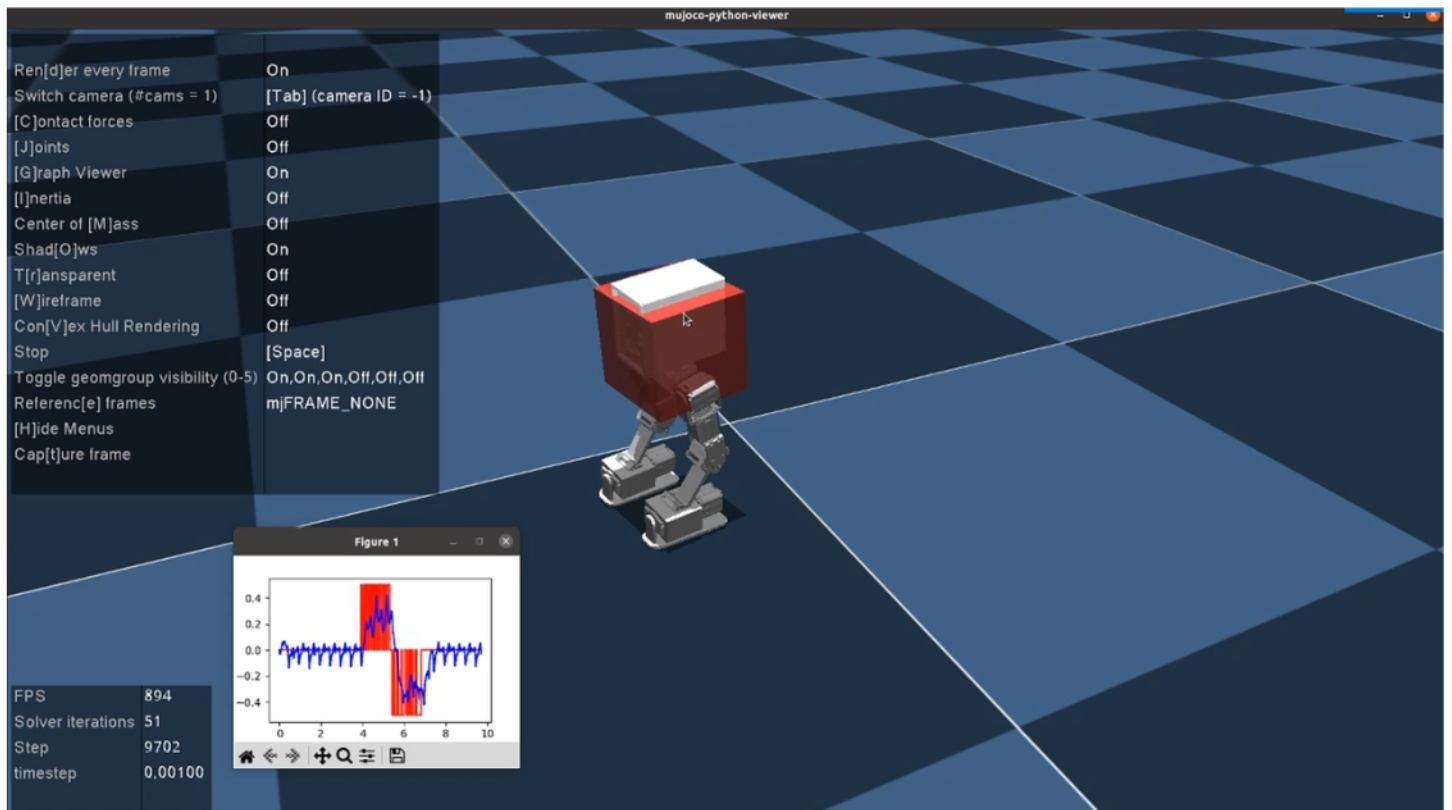
代码块

```
1 python scripts/sim2sim.py --load_model  
  /path/to/logs/Pai_ppo/exported/policies/policy_torch.pt
```

(Note: Replace `/path/to` with the path of the parent folder of `logs`.)



If the operation is successful, the Mujoco scene shown in the figure below will be displayed:

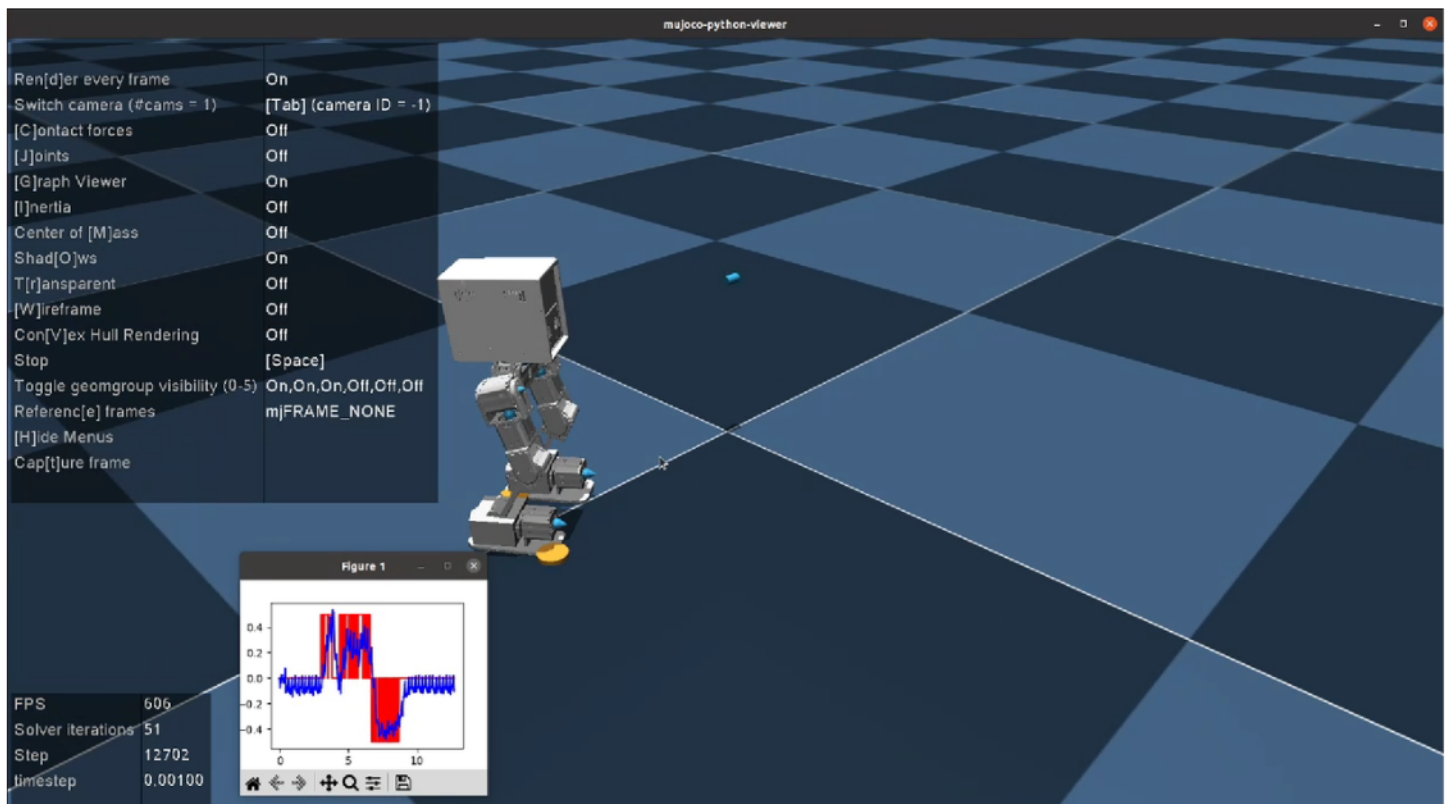


## 16. Run the pre-trained policy\_example provided by High Torque:

代码块

```
1 python scripts/sim2sim.py --load_model
   /path/to/logs/Pai_ppo/exported/policies/policy_example.pt
```

(Note: Replace `/path/to` with the path of the parent folder of `logs`.)





This concludes the basic reinforcement learning baseline tutorial. You can debug your own robot in the simulation environment. For more technical support, please call: +86-18144897433 or add WeChat: GQJD2022.