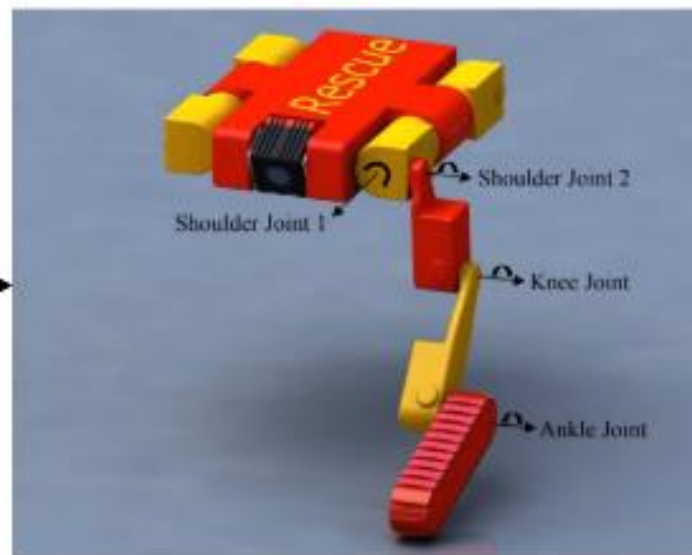


## ENVIRONMENT



## ACTION

joint torsion:  $\tau_{Rji} \dots \tau_{Lji}$   
tracks rotation:  $\tau_{T_{Rji}} \dots \tau_{T_{Lji}}$

## REWARD

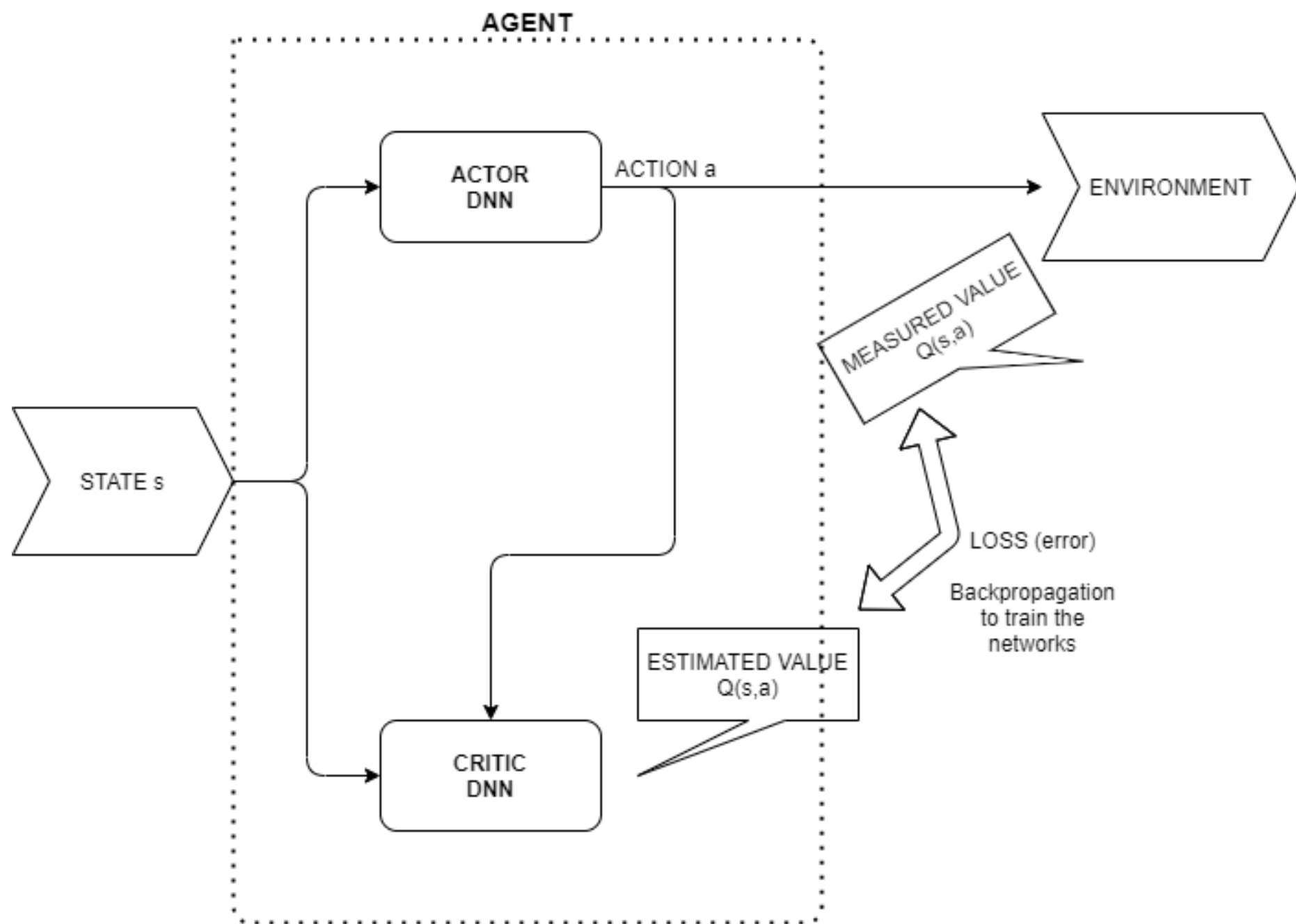
reward function tbd

## STATE

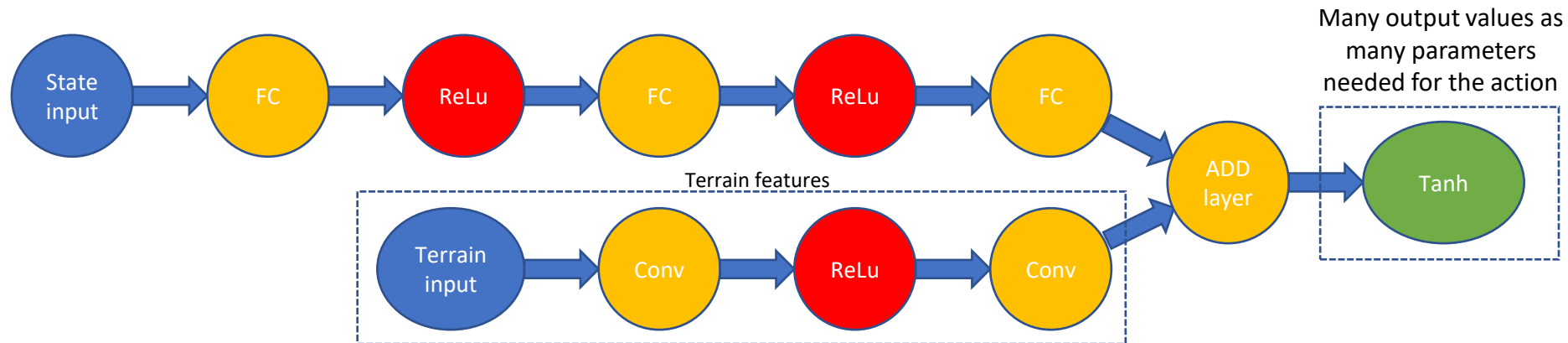
$X, Y, Z$  : translation  
 $\psi, \theta, \phi$  : body rotation  
 $q_{Rji}, q_{Lji}$  : # 16 encoders  
 $F_{Ri}, F_{Li}$  : Normal forces

## AGENT

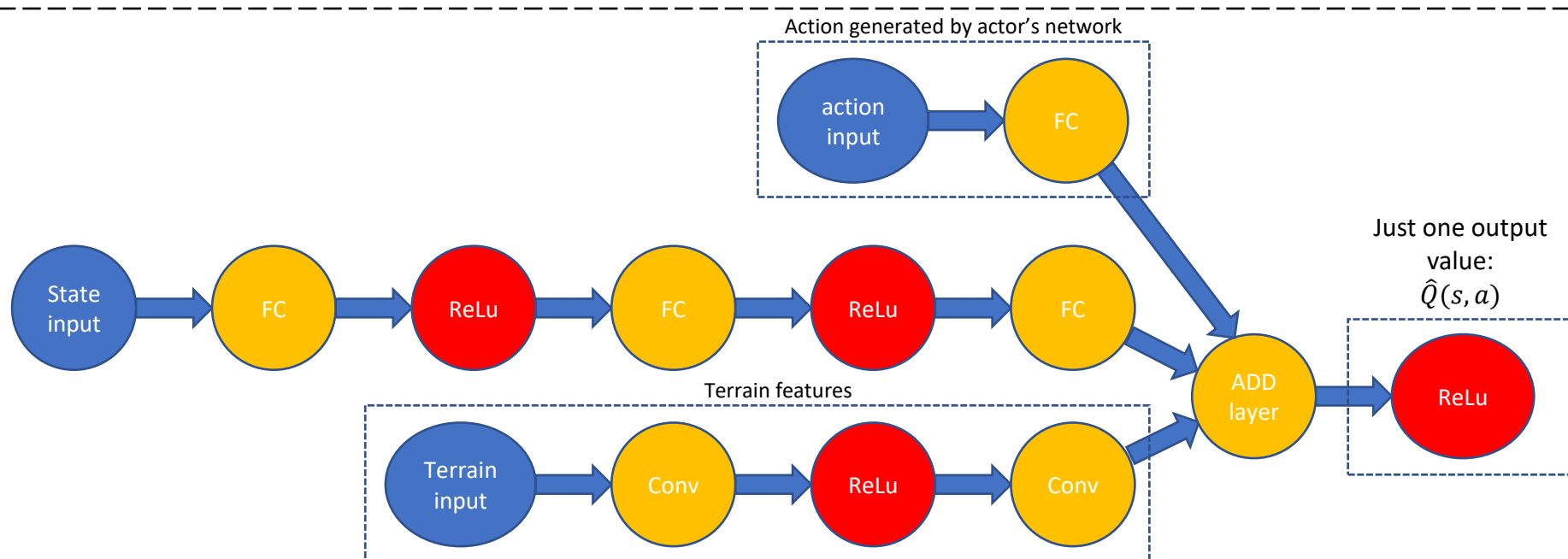
Actor-critic  
Deep Deterministic Policy Gradient  
(DDPG)



## Actor's network



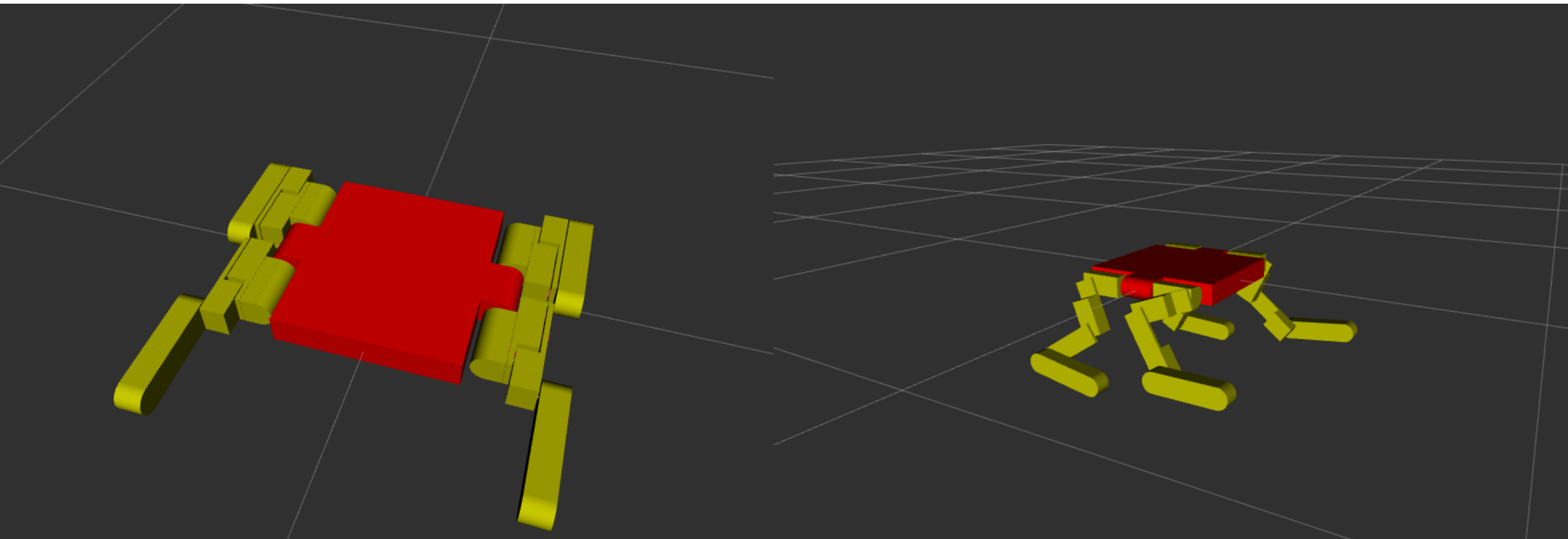
## Critic's network

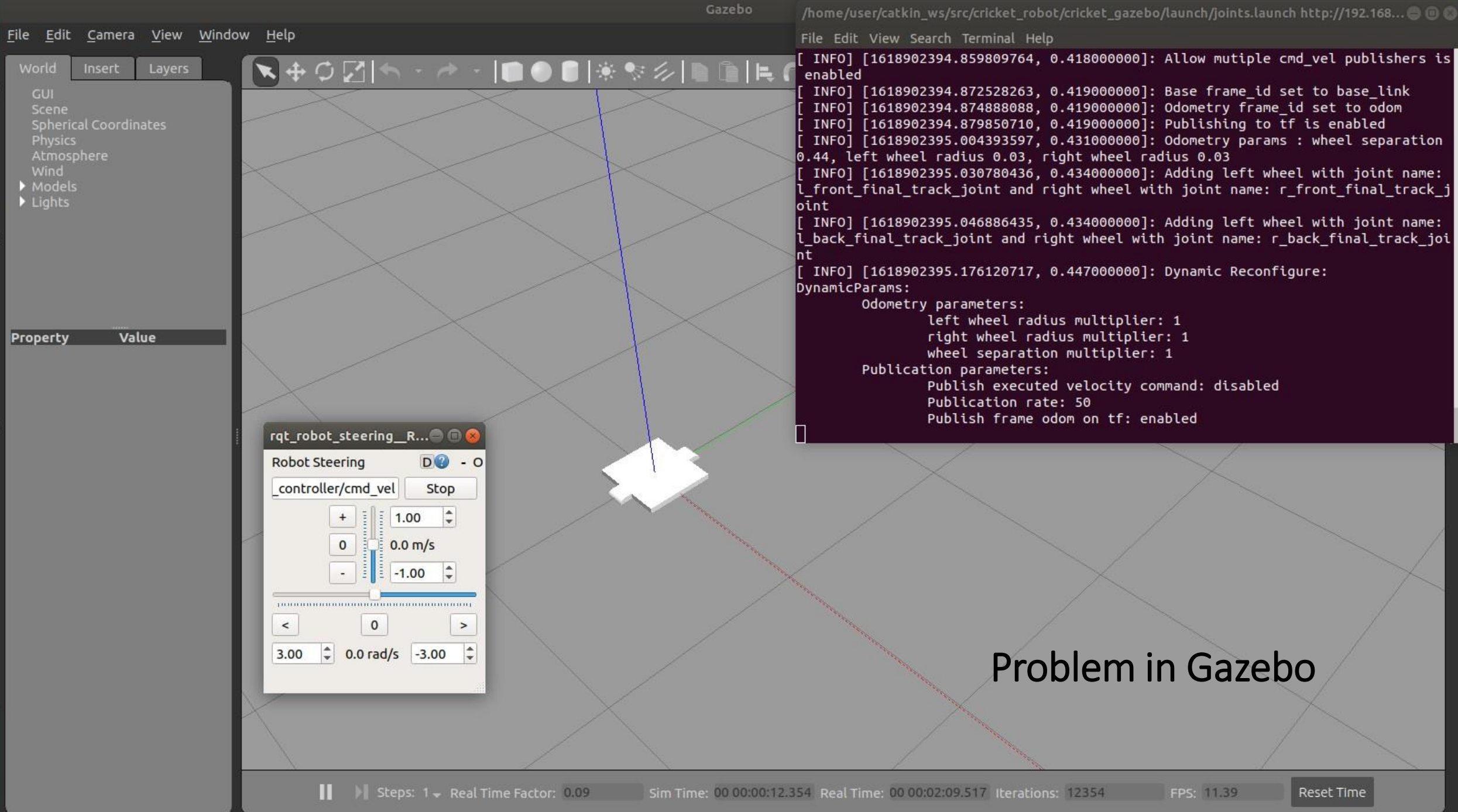


# Reward parameters

Forward scaling → forward reward $X$	A reward for following the right direction to reach the goal
Penalty deviation $Y$	Total deviation penalty
Penalty deviation $Z$	
Torque penalty	Avoid that the robot makes useless movements
Duration reward/penalty	Optimize the time to reach the goal
(optional) energy consumption reward/penalty	Optimize the energy consumption

# URDF model





Problem in Gazebo

It is a reward and not a penalty to avoid the robot decides to encounter stop conditions instead of going on

Avoid straying from decided path

Discount factor  $0 < \gamma \leq 1$

The error  $\varepsilon$  represent the difference between the predicted  $\hat{Q}(s, a)$  and the measured  $Q(s, a)$ .

$$\mathcal{R}_t = v_x + w_t \cdot t - z^2 - w_y y^2 - \sum_i (w^i q_{t-1}^i)^2 - w_\varepsilon \left( \sum_{i=0}^5 \gamma^i \varepsilon_{t-i} \right)^2$$

Forward velocity

Distance from basic height

Minimize the encoders effort (considering the previous values of  $q$ )

This gives the robot a way to understand the error term and develop a policy accordingly

It is a reward and not a penalty to avoid the robot decides to encounter stop conditions instead of going on

Avoid straying from decided path

Discount factor  
 $0 < \gamma \leq 1$

The error  $\varepsilon$  represent the difference between the predicted  $\hat{Q}(s, a)$  and the measured  $Q(s, a)$ .

$$\mathcal{R}_t = \underbrace{v_x}_{\text{Forward velocity}} + \underbrace{0.0625 \cdot t}_{\text{It is a reward and not a penalty to avoid the robot decides to encounter stop conditions instead of going on}} - \underbrace{z^2}_{\text{Distance from basic height}} - \underbrace{10y^2}_{\text{Avoid straying from decided path}} - \underbrace{0.02 \sum_i (q_{t-1}^i)^2}_{\text{Minimize the encoders effort (considering the previous values of } q)} - \underbrace{100 \left( \sum_{i=0}^5 \gamma^i \varepsilon_{t-i} \right)^2}_{\substack{\text{This gives the robot a way to understand the error term and develop a policy accordingly} \\ \text{The error } \varepsilon \text{ represent the difference between the predicted } \hat{Q}(s, a) \text{ and the measured } Q(s, a).}}$$



Number of states (input layer dimension for the actor's neural network)

Robot position	$X, Y, Z$	3
Robot velocity	$v_x, v_y, v_z$	3
Robot angle	$\phi, \theta, \psi$	3
Robot angle rate	$\omega_\phi, \omega_\theta, \omega_\psi$	3
4 legs	four joints per leg	16
4 tracks	4 track rotation values	4
Contact forces	6 normal forces per track	24
Computed errors	$\varepsilon_t$	5
Total		61

$$\mathcal{R}_t = - \sum_{i=1}^{no\ tracks} \left[ \alpha_i (q_{t-1}^i)^2 + \beta_i (q_{t-1}^i)^2 + \kappa_\tau^i |\tau_{t-1}^i| + w_q^i (\Delta q_t^i)^2 \right] - w_\epsilon \left( \sum_{i=0}^5 \gamma^i \epsilon_{t-i} \right)^2 - w_t t - \Delta X - \Delta Y - \Delta Z$$

$\beta_i = 0$  or  $1$   
 Penalty if the robot touches the environment

$\alpha_i = 0$  or  $1$   
 Penalty if the robot touches itself

$\kappa_\tau^i = -1$  or  $1$   
 reward/penalty based on the direction of the last rotation

Difference from final position

Discount factor  
 $0 < \gamma \leq 1$

The error  $\epsilon$  represent the difference between the **predicted**  $\hat{Q}(s, a)$  and the **measured**  $Q(s, a)$ .

This gives the robot a way to understand the error term and develop a policy accordingly

Time passed

Difference from final optimal position

$$x_{norm} = \frac{x - x_{min}}{x_{max} - x_{min}}$$

$$\mathcal{R}_t = - \sum_{i=1}^{no\ tracks} \left[ \alpha_i (q_{t-1}^i)^2 + \beta_i (q_{t-1}^i)^2 + \kappa_\tau^i |\tau_{t-1}^i| + w_q^i (\Delta q_t^i)^2 \right] - w_\varepsilon \left( \sum_{i=0}^5 \gamma^i \varepsilon_{t-i} \right)^2 - w_t t - w_X (\Delta X)^2 - w_Y (\Delta Y)^2 - w_Z (\Delta Z)^2 - w_\psi (\Delta \psi)^2 - w_\theta (\Delta \theta)^2 - w_\phi (\Delta \phi)^2$$

The error  $\varepsilon$  represent the difference between the **predicted**  $\hat{Q}(s, a)$  and the **measured**  $Q(s, a)$ .

$\beta_i = 0$  or  $1$   
Penalty if the robot touches the environment

Difference with the joints final position

Discount factor  $0 < \gamma \leq 1$

$\alpha_i = 0$  or  $1$   
Penalty if the robot touches itself

$\kappa_\tau^i = -1$  or  $1$   
reward/penalty based on the direction of the last rotation

This gives the robot a way to understand the error term and develop a policy accordingly

Difference with the robot's final centre of mass position

Time passed

Difference with the robot's final body rotation

$$\mathcal{R}_t = - \sum_{i=1}^{no\ tracks} \left[ \alpha_i (q_{t-1}^i)^2 + \beta_i (q_{t-1}^i)^2 + \kappa_\tau^i |\tau_{t-1}^i| + w_q^i (\Delta q_t^i)^2 \right] - w_\varepsilon \left( \sum_{i=0}^5 \gamma^i \varepsilon_{t-i} \right)^2 - w_t t - w_X (\Delta X)^2 - w_Y (\Delta Y)^2 - w_Z (\Delta Z)^2 - w_\psi (\Delta \psi)^2 - w_\theta (\Delta \theta)^2 - w_\phi (\Delta \phi)^2$$

# Deep Deterministic Policy Gradient

[1] Lillicrap, Timothy P., et al. "Continuous control with deep reinforcement learning" (2015)

[2] Zhipeng L., Hao C., et al. "Adversarial Deep Reinforcement Learning in Portfolio Management" (2018)

```
1: Initialize actor and critic networks  $\mathcal{A}(s, e|\theta^{\mathcal{A}})$  and  $\mathcal{C}(s, a, e|\theta^{\mathcal{C}})$  with random weights  $\theta^{\mathcal{A}}$  and  $\theta^{\mathcal{C}}$ 
2: Initialize target network  $\mathcal{A}'$  and  $\mathcal{C}'$  with random weights  $\theta^{\mathcal{A}'} \leftarrow \theta^{\mathcal{A}}$  and  $\theta^{\mathcal{C}'} \leftarrow \theta^{\mathcal{C}}$ 
3: Initialize replay buffer  $\mathcal{R}$ 
4: for  $n=1$  to  $N_{episodes}$  do:

5:     Procedure RobotInit
6:     Generate a random stable starting state  $s$  in the unstructured environment
7:     end procedure RobotInit

8:     Procedure Simulation
9:     for  $t=0$  to  $T$  do:
10:         generate action  $a_t \leftarrow \mathcal{A}(s_t, e|\theta^{\mathcal{A}})$ 
11:         Execute action  $a_t$ : observe reward  $r_t$  and new state  $s_{t+1}$ 
12:         Store transition  $(s_t, a_t, r_t, s_{t+1})$  in  $\mathcal{R}$ 
13:         Sample a random minibatch of  $N$  transitions  $(s, a, r, s')$  from  $\mathcal{R}$ 
14:         for  $(s_i, a_i, r_i, s_{i+1})$  in minibatch do:
15:             set  $y_i = r_i + \gamma \cdot \mathcal{C}'(s_{i+1}, \mathcal{A}'(s_{i+1}, e|\theta^{\mathcal{A}'}), e|\theta^{\mathcal{C}'})$ 
16:         end for
17:         Update critic by minimizing the loss  $L = \frac{1}{N} \sum_i (y_i - \mathcal{C}(s_i, a_i, e|\theta^{\mathcal{C}}))^2$ 
18:         Update actor policy by policy gradient
19:         Update target networks  $\mathcal{A}'$  and  $\mathcal{C}'$ 
20:         checkStopCondition()
21:     end for
22:     end procedure Simulation
23: end for
```