

Implementation of Parking Using Reinforcement Learning

2조 : 김현우
김범진
김창기
안종원

01 Project Outline

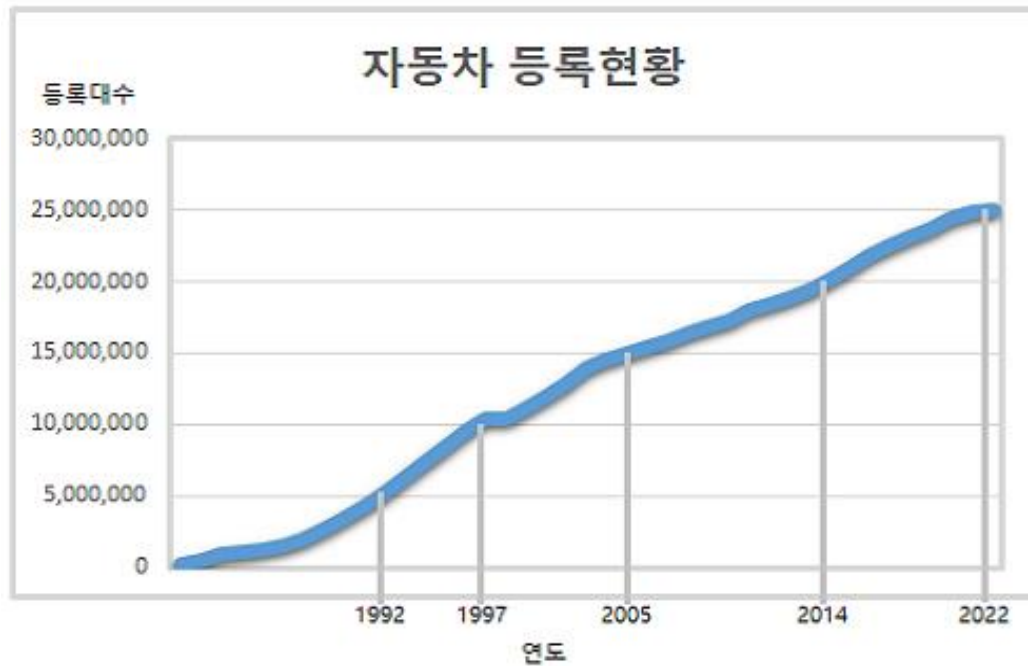
02 Project Implementation

03 Project Result

1

Project Outline

Project Outline



상반기 국내 자동차 시장 키워드 (단위: %)

※ 출처: 한국자동차산업협회(KAMA)

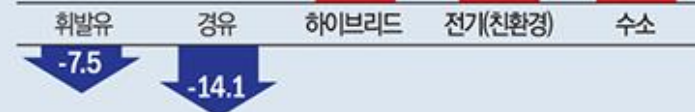
대형화

※ 크기별 전년 대비 증감률



전동화

※ 동력원별 신규등록
전년 대비 증감률



고급화

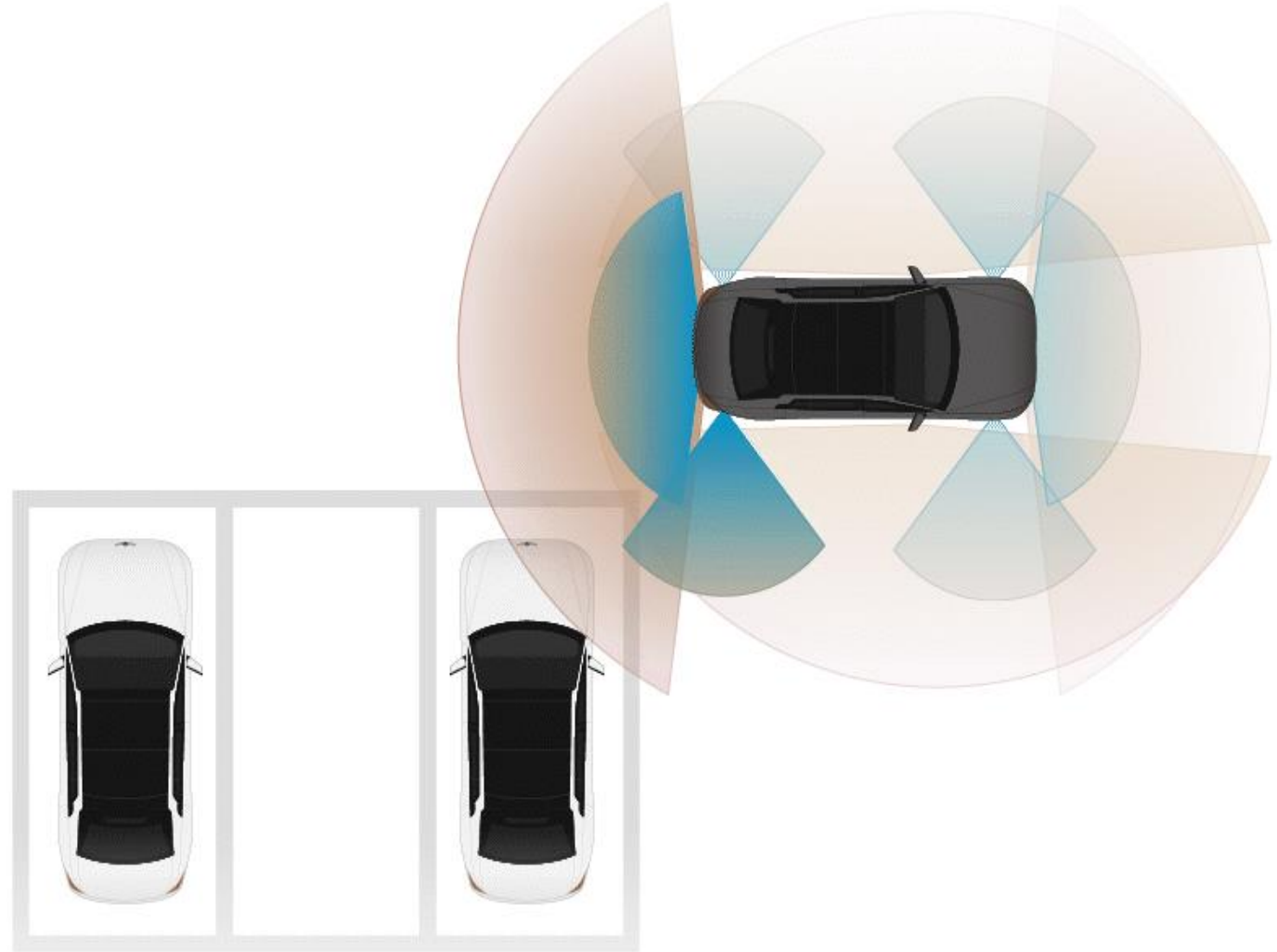
※ 전년 대비 증감률



Project Outline



Project Outline

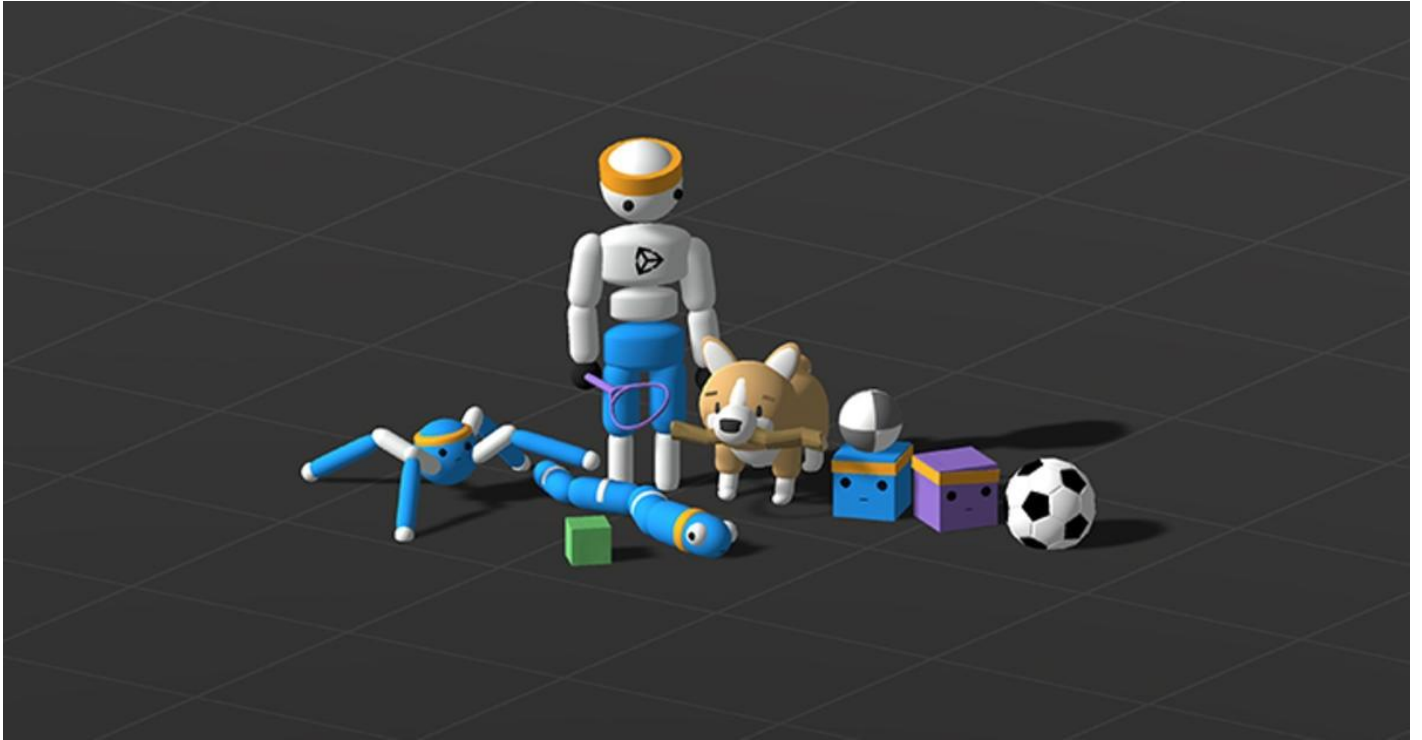


2

Project Implementation

Project Implementation

Use Unity ML-Agent to implement



Project Implementation



Project Implementation

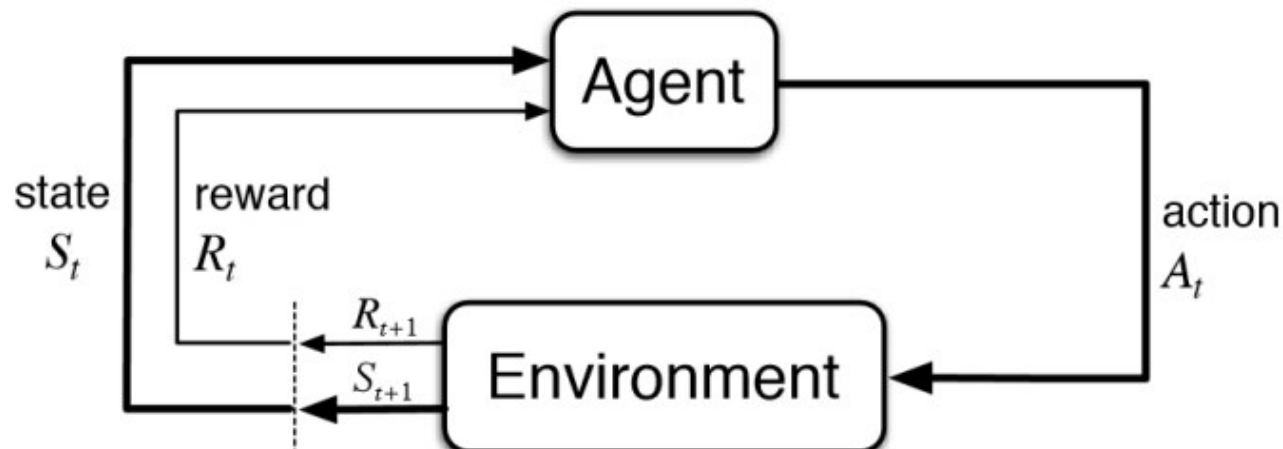
Use Unity Asset Store's Car asset and town asset.



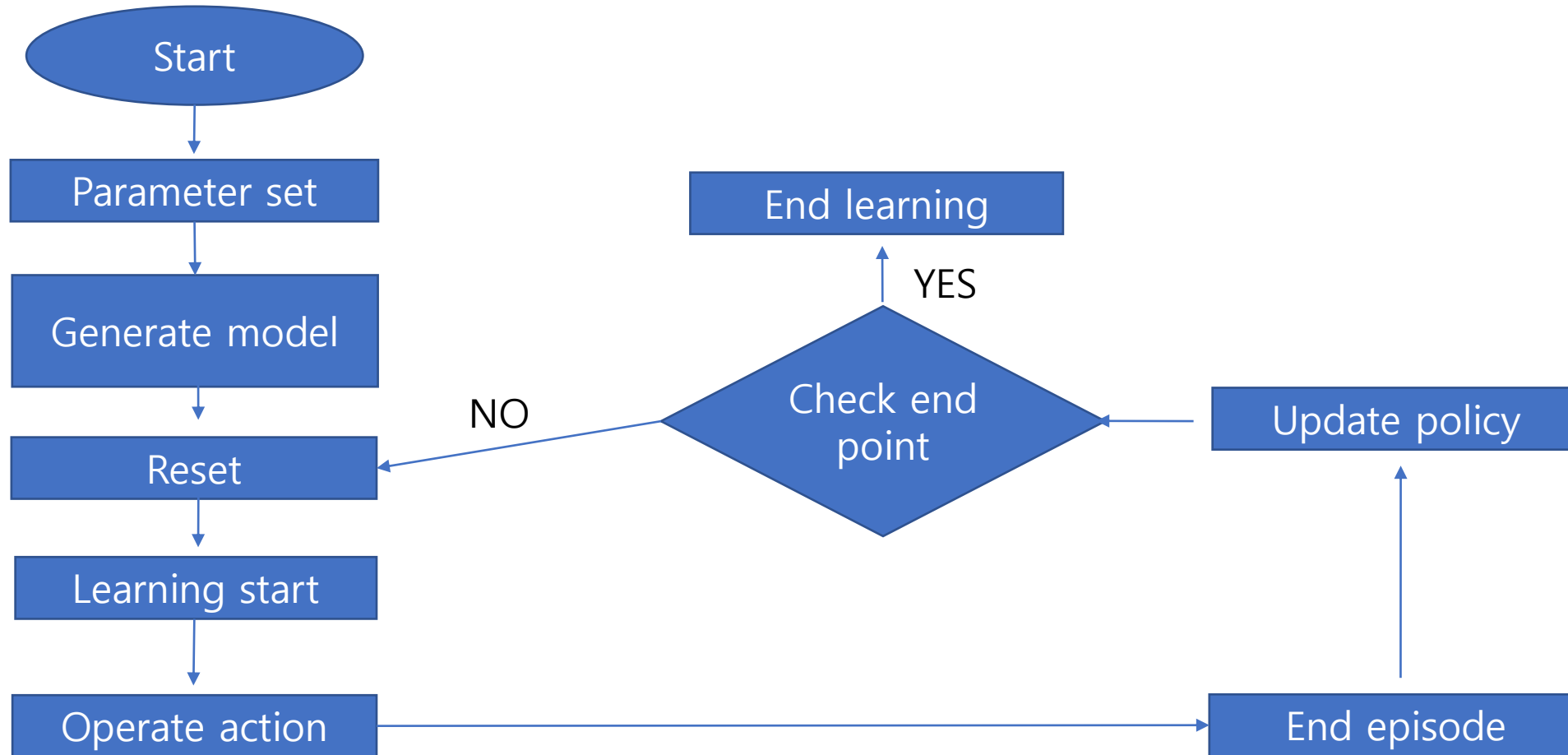
Project Implementation

Reinforcement Learning?

- Agent receives a State from the environment to determine an action
- Environment takes action from the agent, performs, and then outputs the state and the return
- Learn an action that maximizes the expected value by the interaction between the environment and the agent.



Project Implementation



Project Implementation

Reinforcement method: PPO(Proximal Policy Optimization)

Algorithm 1 PPO-Clip

- 1: Input: initial policy parameters θ_0 , initial value function parameters ϕ_0
- 2: **for** $k = 0, 1, 2, \dots$ **do**
- 3: Collect set of trajectories $\mathcal{D}_k = \{\tau_i\}$ by running policy $\pi_k = \pi(\theta_k)$ in the environment.
- 4: Compute rewards-to-go \hat{R}_t .
- 5: Compute advantage estimates, \hat{A}_t (using any method of advantage estimation) based on the current value function V_{ϕ_k} .
- 6: Update the policy by maximizing the PPO-Clip objective:

$$\theta_{k+1} = \arg \max_{\theta} \frac{1}{|\mathcal{D}_k|T} \sum_{\tau \in \mathcal{D}_k} \sum_{t=0}^T \min \left(\frac{\pi_{\theta}(a_t|s_t)}{\pi_{\theta_k}(a_t|s_t)} A^{\pi_{\theta_k}}(s_t, a_t), \quad g(\epsilon, A^{\pi_{\theta_k}}(s_t, a_t)) \right),$$

typically via stochastic gradient ascent with Adam.

- 7: Fit value function by regression on mean-squared error:

$$\phi_{k+1} = \arg \min_{\phi} \frac{1}{|\mathcal{D}_k|T} \sum_{\tau \in \mathcal{D}_k} \sum_{t=0}^T \left(V_{\phi}(s_t) - \hat{R}_t \right)^2,$$

typically via some gradient descent algorithm.

- 8: **end for**
-

Project Implementation

Reinforcement method : PPO(Proximal Policy Optimization)

```
1 behaviors:
2   CarBehaviour:
3     trainer_type: ppo
4     hyperparameters:
5       batch_size: 1024
6       buffer_size: 5120
7       learning_rate: 0.00025
8       beta: 0.0020
9       epsilon: 0.15
10      lambda: 0.95
11      num_epoch: 5
12      learning_rate_schedule: linear
13
14    network_settings:
15      normalize: true
16      hidden_units: 264
17      num_layers: 3
18
19    reward_signals:
20      extrinsic:
21        gamma: 0.95
22        strength: 0.99
23
24    keep_checkpoints: 15
25    checkpoint_interval: 100000
26    time_horizon: 264
27    max_steps: 100000000
28    summary_freq: 100000
```

Project Implementation

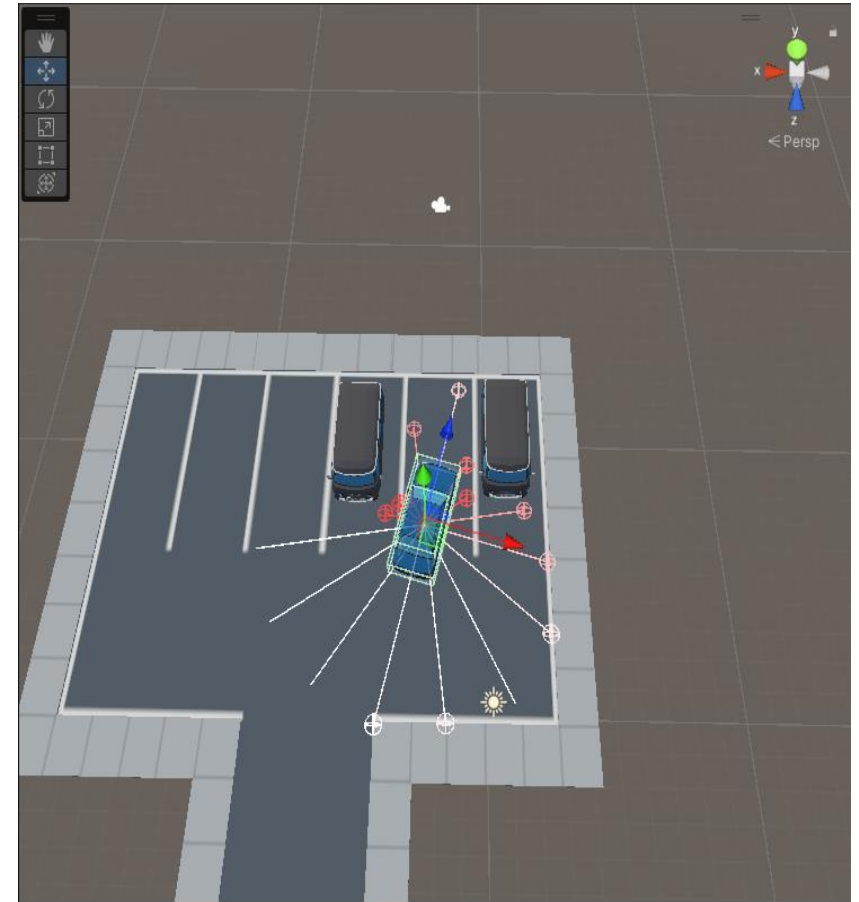
◆ Action : steer, forward, backward



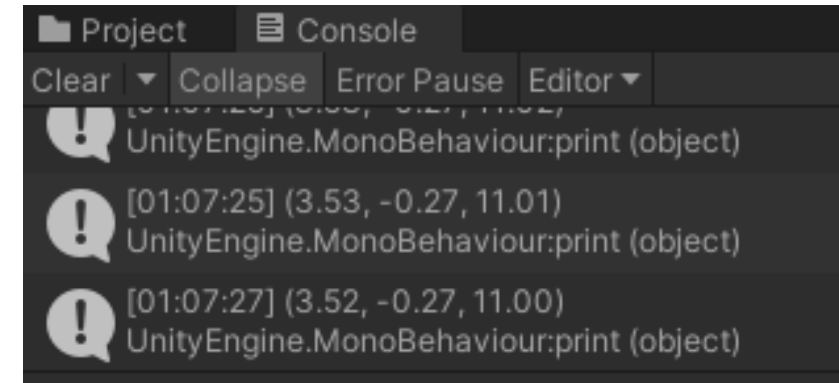
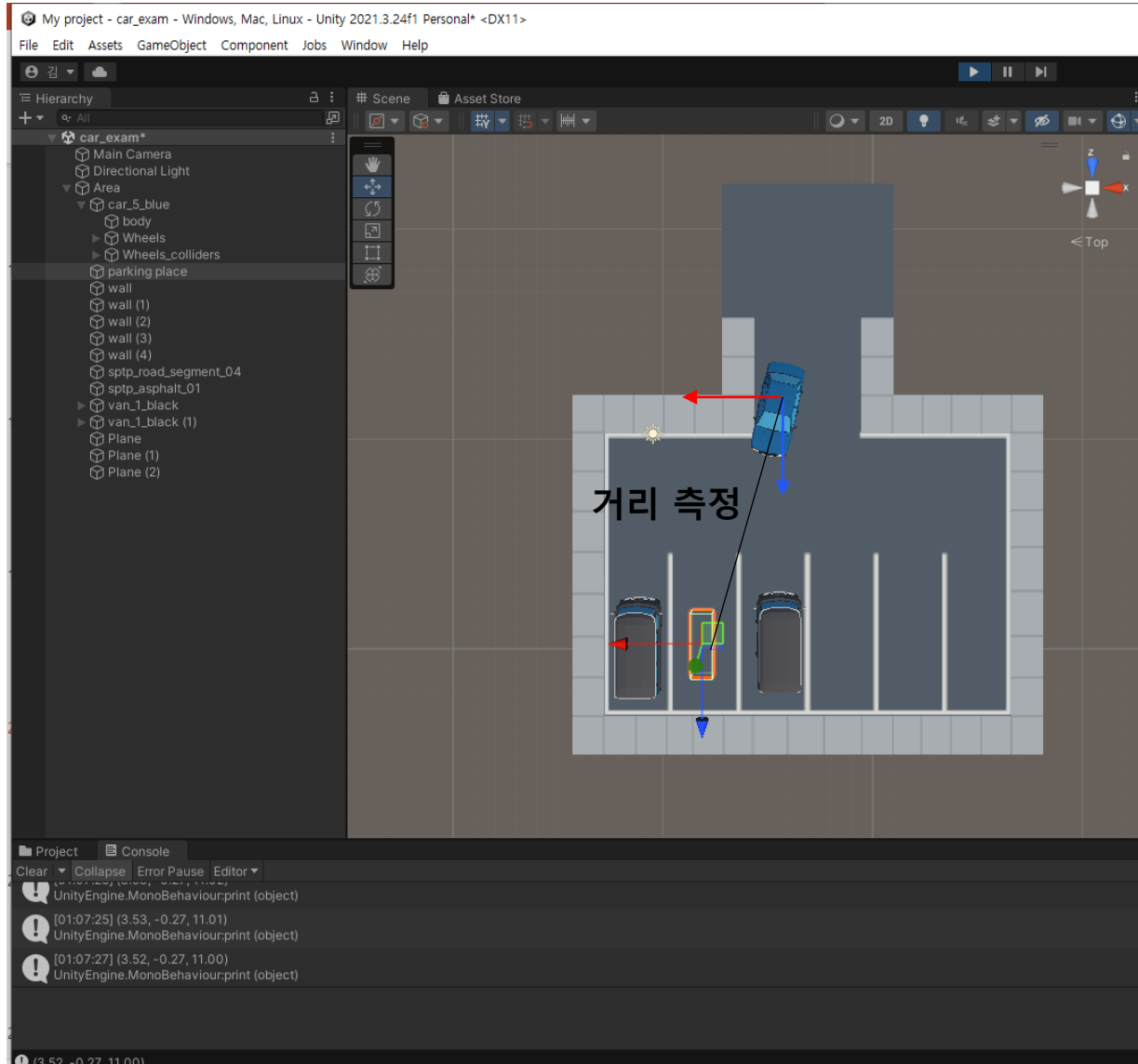
Project Implementation

◆ State :

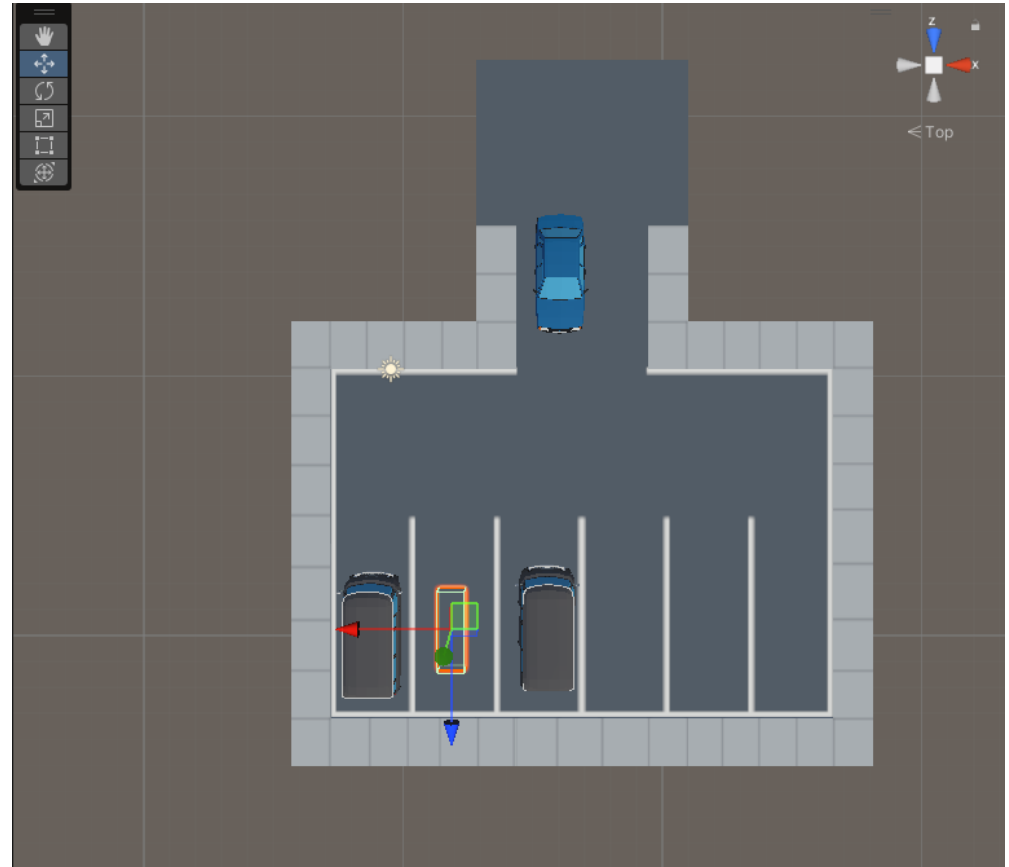
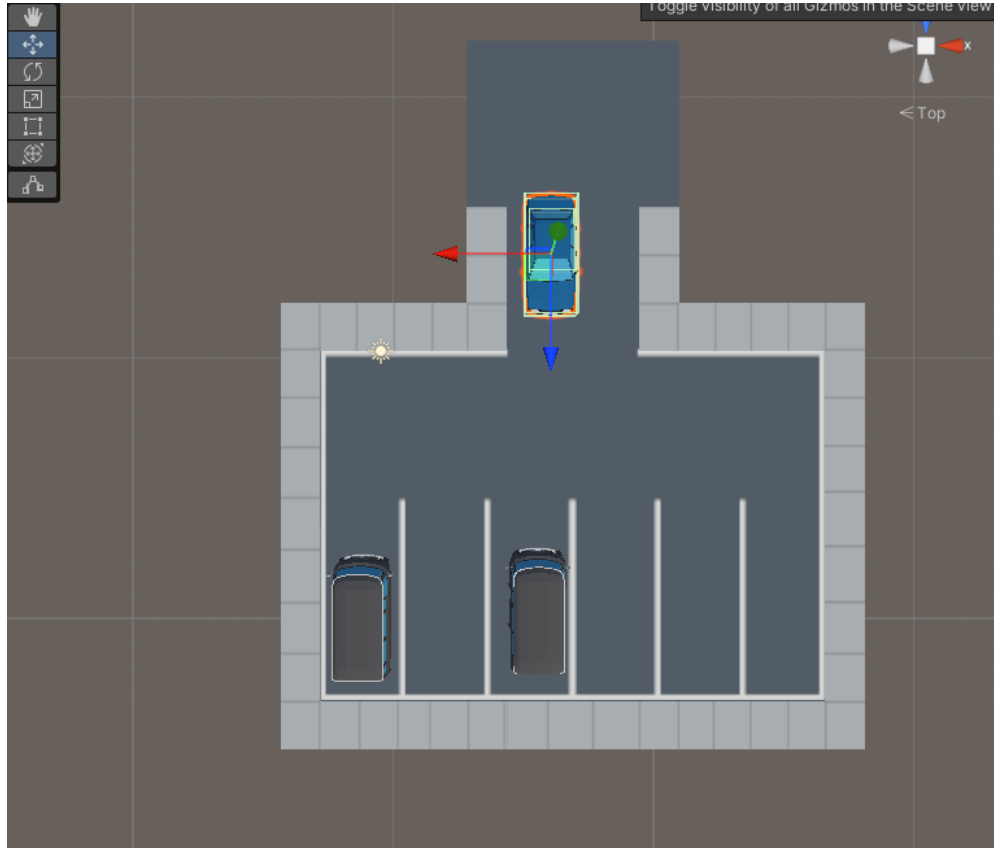
1. Agent's Position – Parking place Position (x,y,z)
 2. Agent's angle – Parking place angle
 3. Agent's speed
- +RayPerceptionSensor 3D sensing



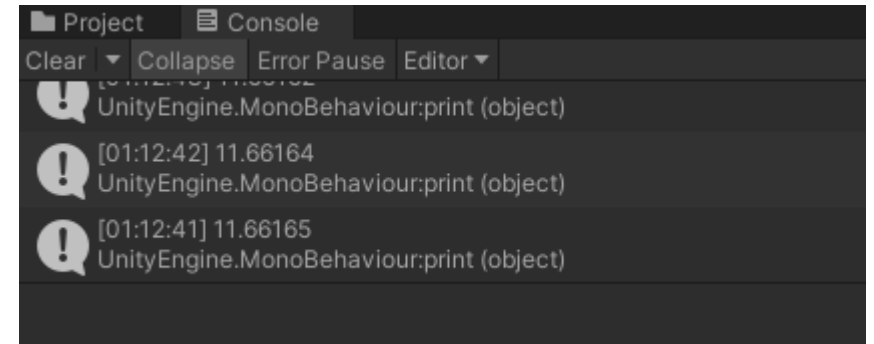
Project Implementation



Project Implementation



Project Implementation



Project Implementation

◆ Reward :

1. Parking success+10reward
2. When Collision –reward depend on Agent speed
3. Agent position – Parking place position
4. Agent angle – Parking place angle

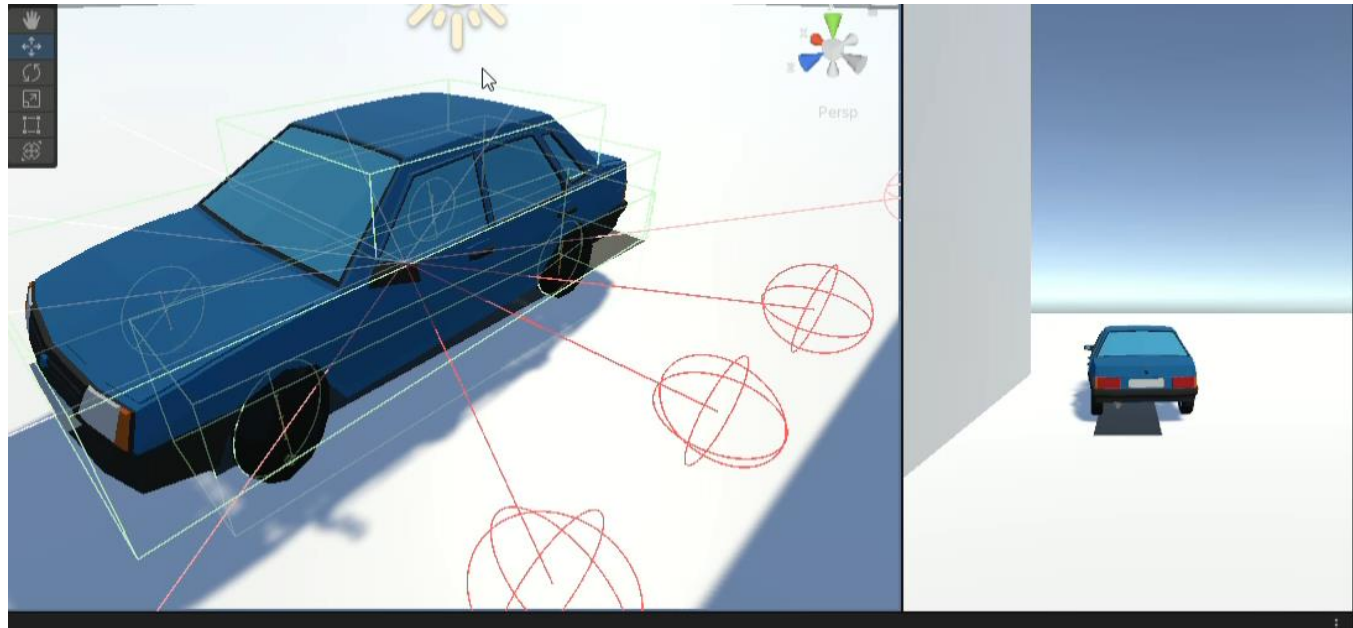


3

Project Result

Project Result

Empty Place



Project Result

Only Parking line

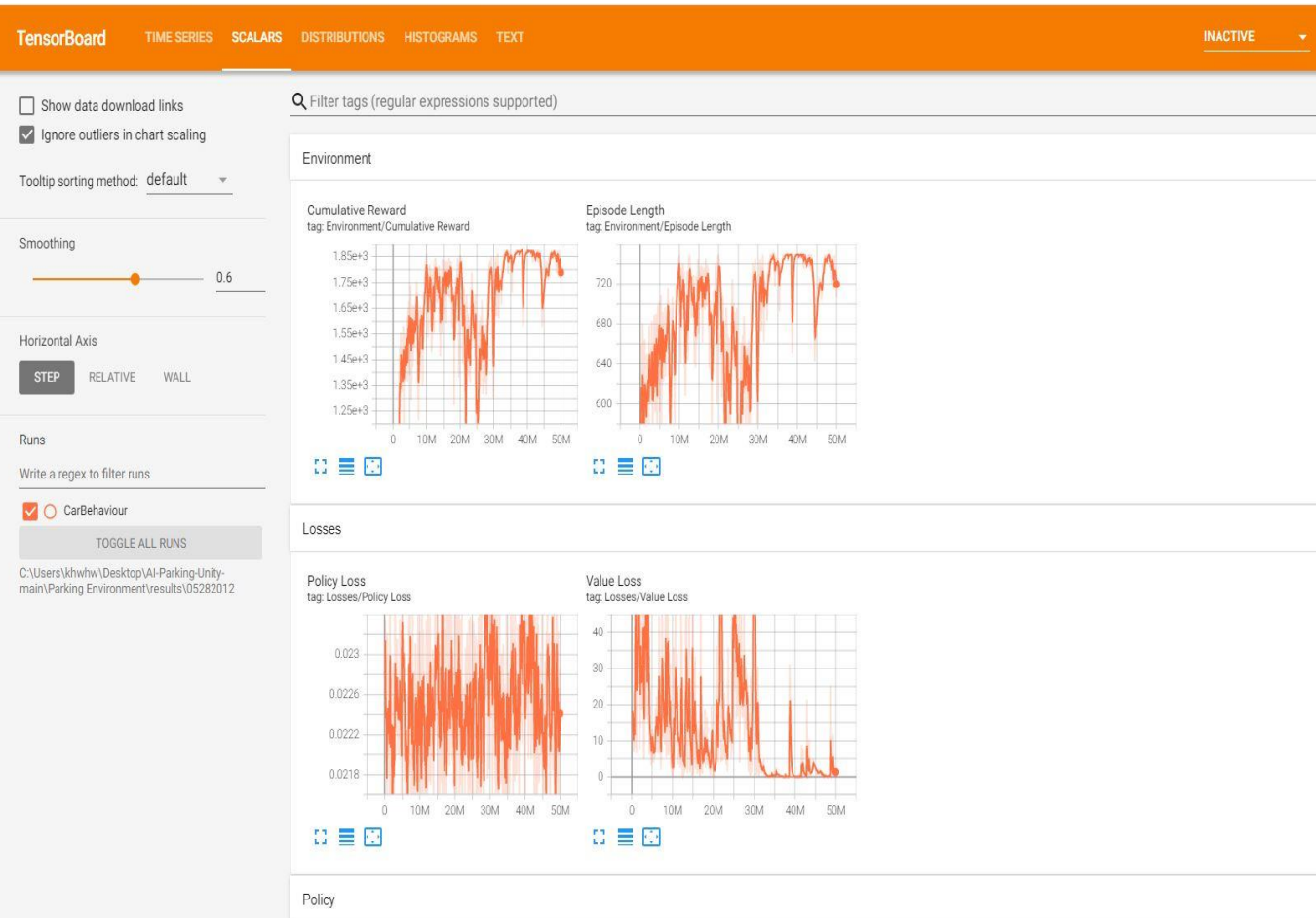
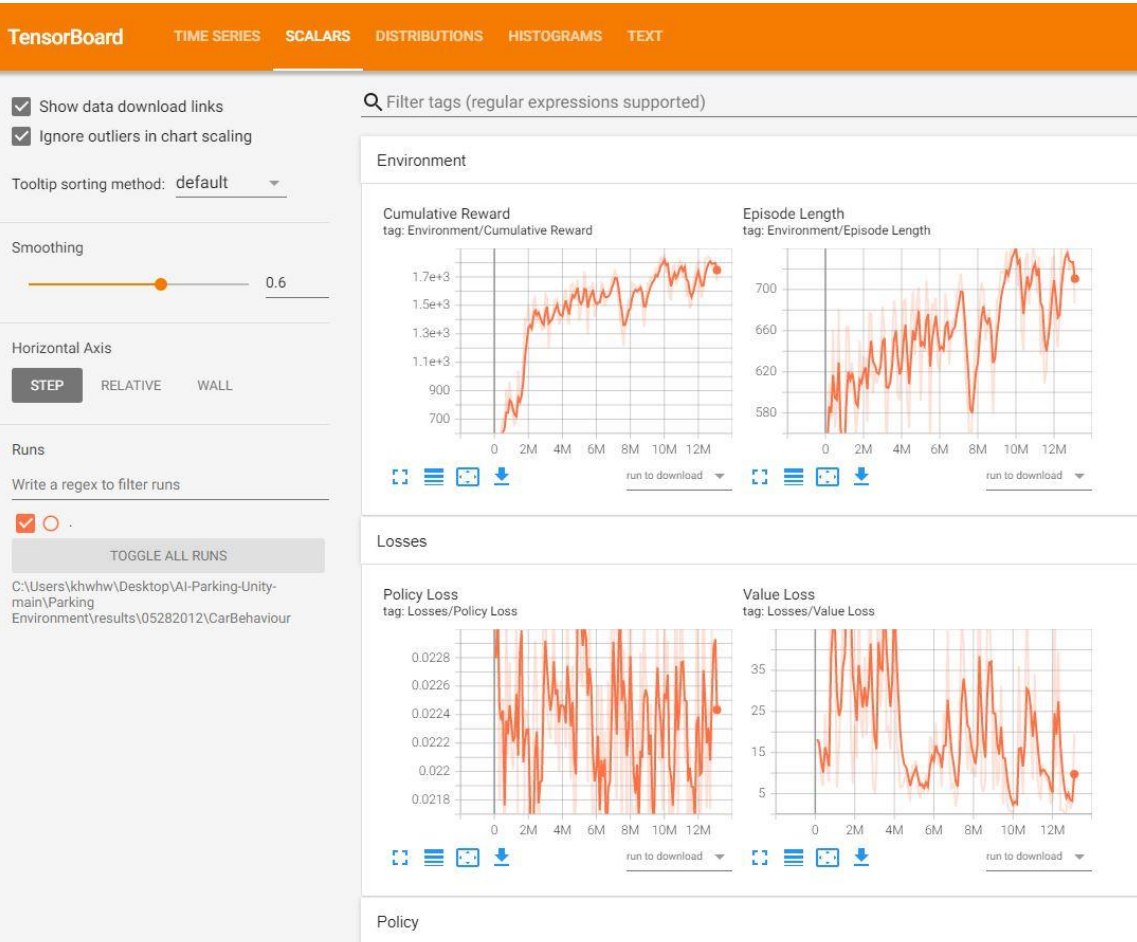


Project Result

In constrained situation(Parking line, other Cars)



Project Result

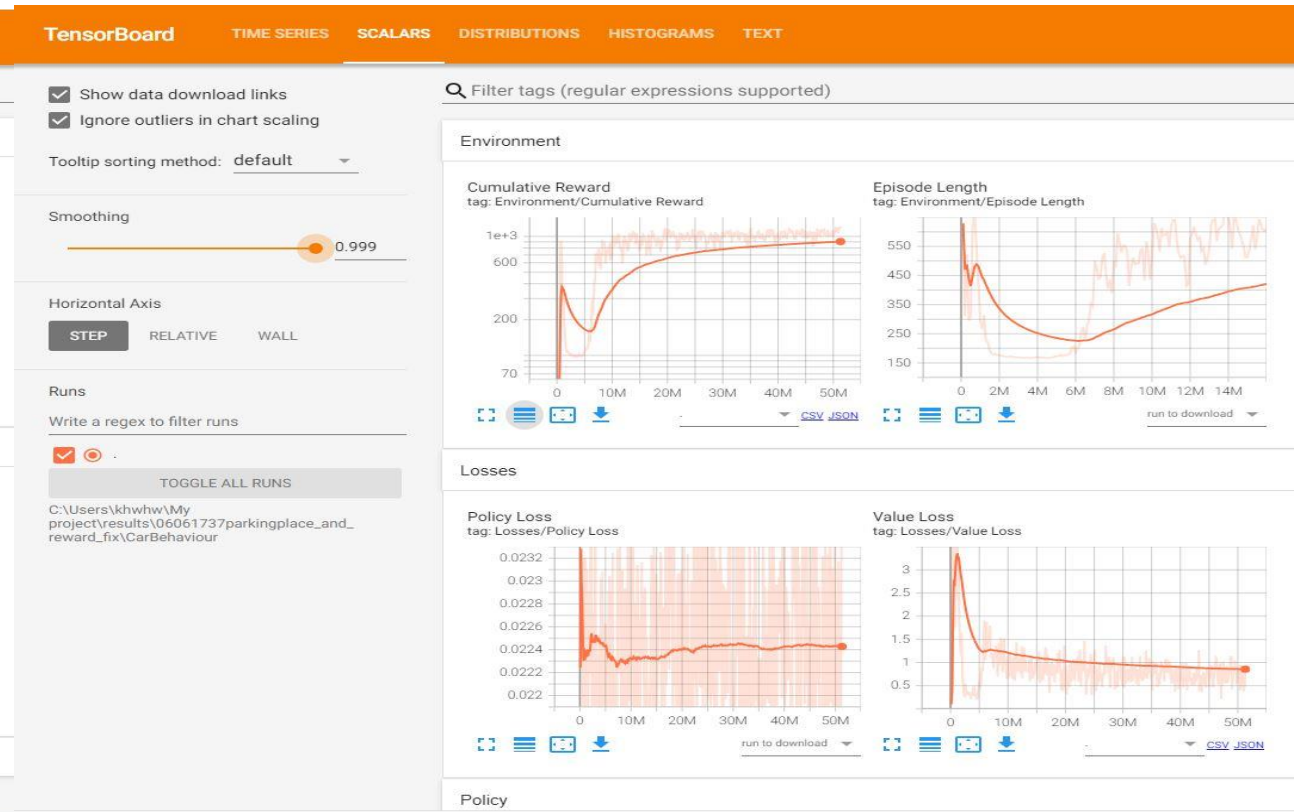
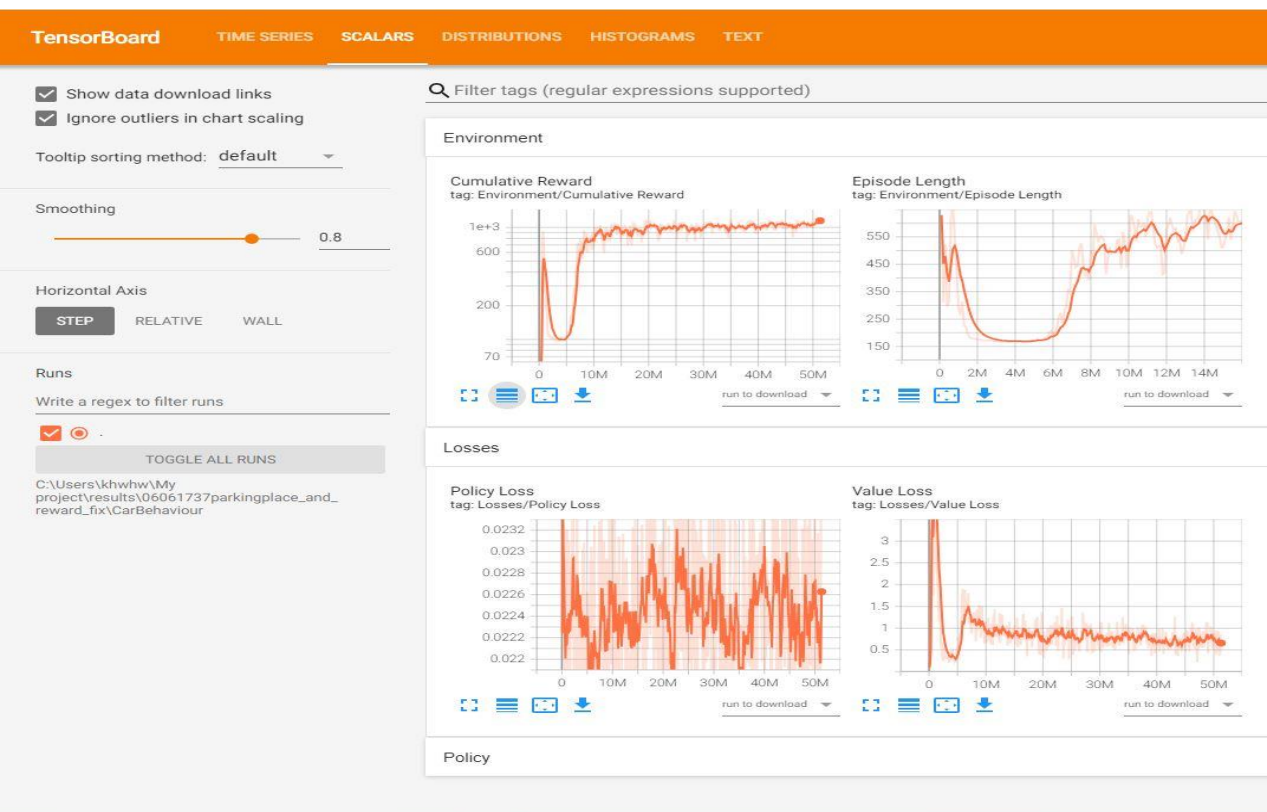


Project Result

Make reward more stably



Project Result



Thank you!
