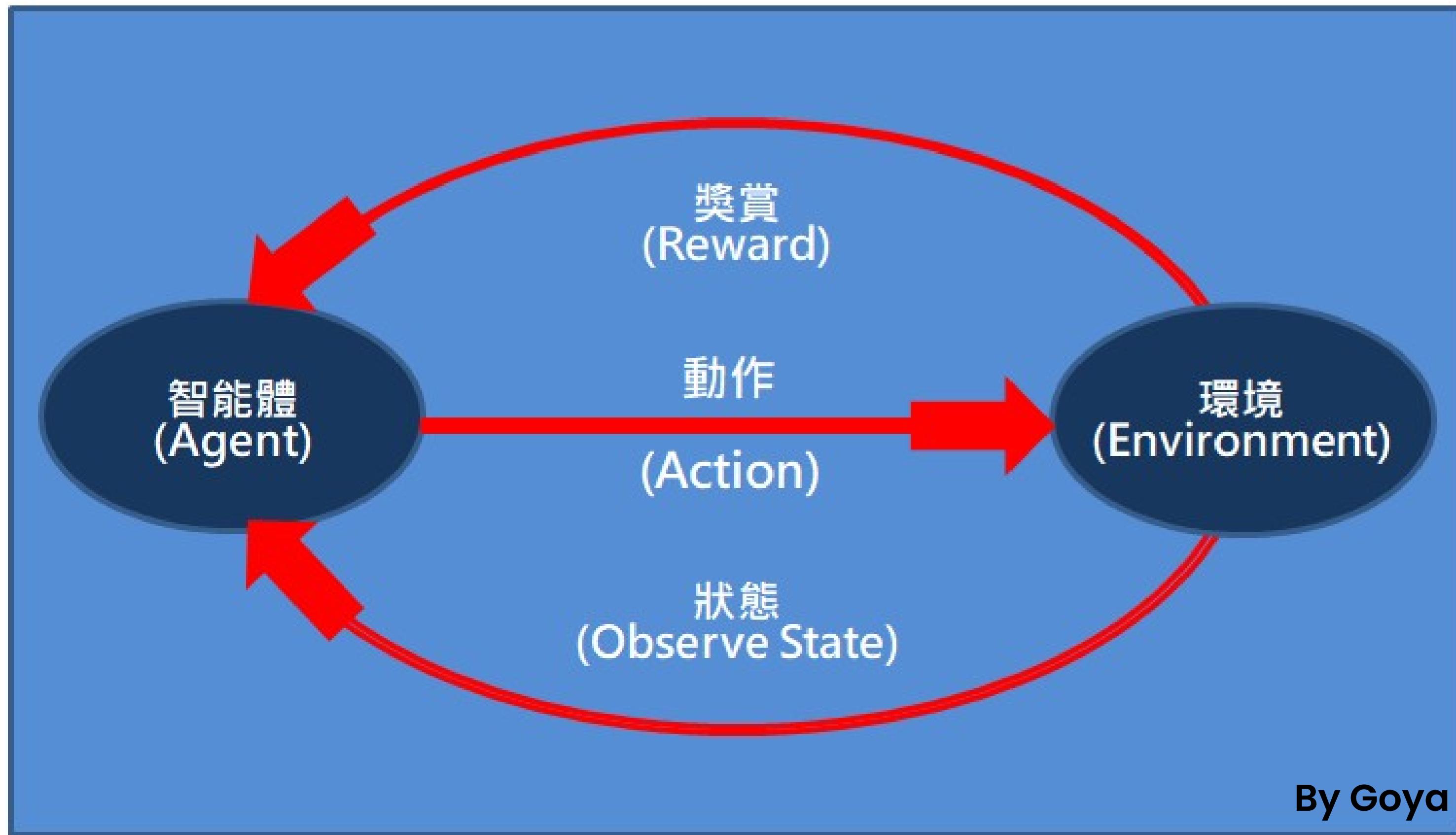


# 深度學習期末報告

Use Unity ML-Agent  
to make Ai parking

組員:周翰文、李沐風

# 強化學習簡介



# 算法:PPO(近便政策最佳化)

## 策略梯度方法

數據效率低

策略更新穩定

# 使用工具

# 工具列表



# 前置環境

(設定程式、設定場景、設定訓練陣列、設定傳感器)



# 程式

# Behavior Parameters

Behavior Name: CarBehaviour

Vector Observation: Space Size: 3, Stacked Vectors: 1

Actions: Continuous Actions: 3, Discrete Branches: 0

Model: Inference Device: GPU, Behavior Type: Default, Team Id: 0, Use Child Sensors: checked, Observable Attribute Handler: Ignore

# Demonstration Recorder

Record: checked, Num Steps To Record: 0, Demonstration Name: 90gays, Demonstration Directory: Demos

# Car Agent (Script)

Max Step: 750

Script: CarAgent

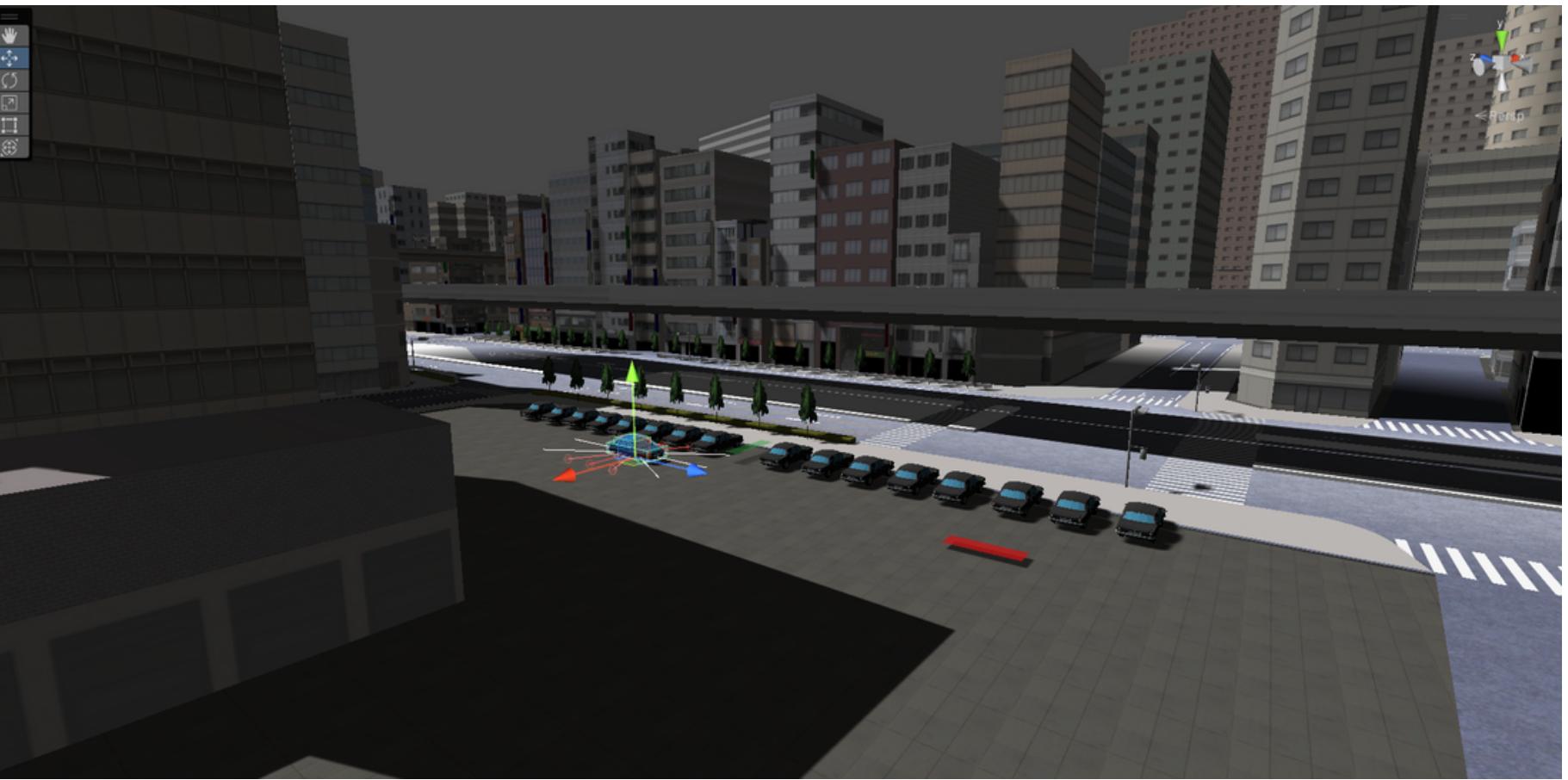
Spawn Radius X: 6, Spawn Radius Z: 0.25, Env Radius X: 7.5, Env Radius Z: 7.5, In Target Multiplier: 1.5, Target: 0

Car Cameras: List is Empty

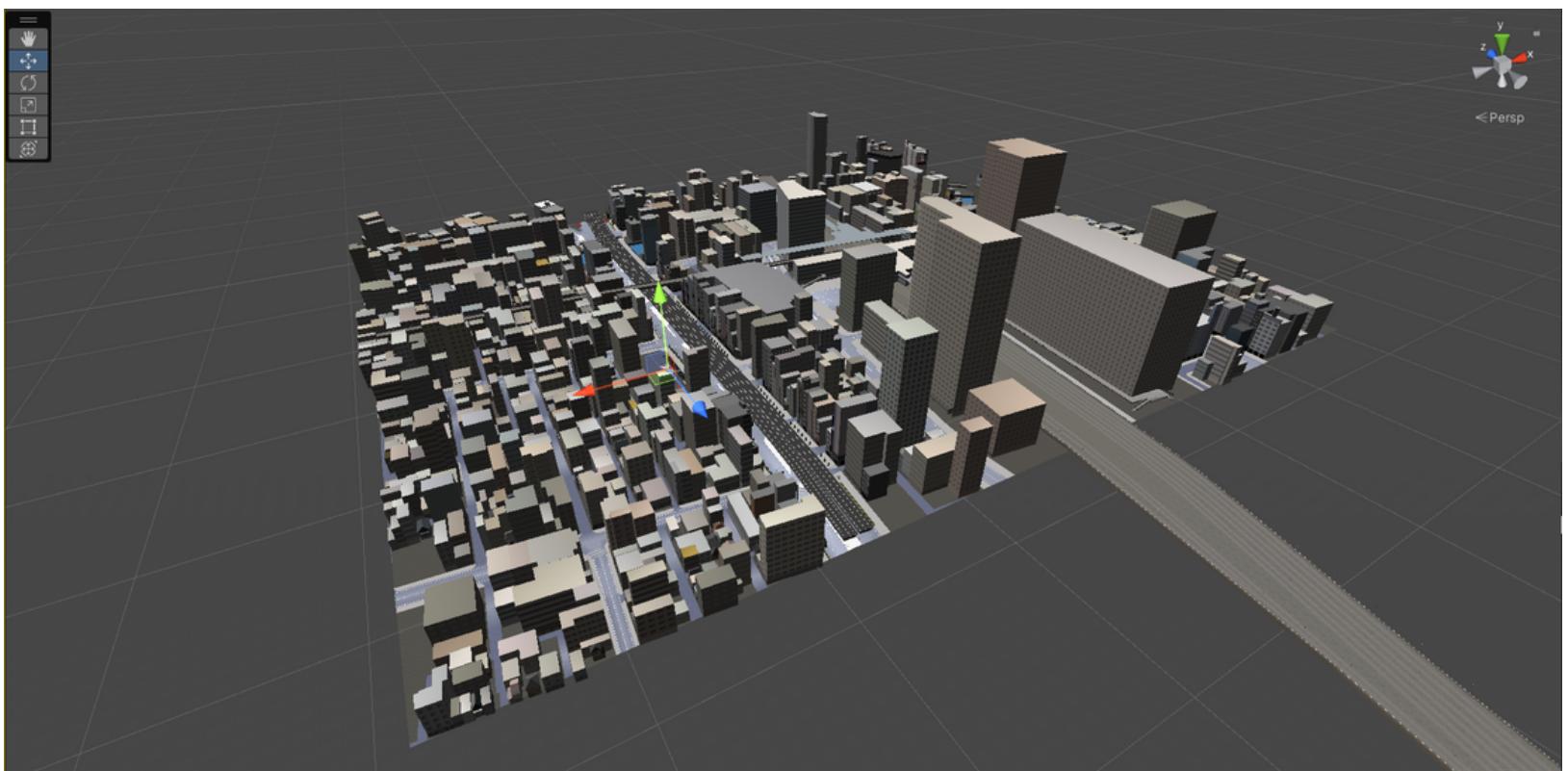
Camera Width: 200, Camera Height: 200, Camera Gray Scale: false, Sensor Compression Type: PNG, Find Parking Spot: false

```
Assembly-CSharp
1 using System.Collections;
2 using System.Collections.Generic;
3 using System.Linq;
4 using UnityEngine;
5 
6 using Unity.MLAgents;
7 using Unity.MLAgents.Sensors;
8 using Unity.MLAgents.Actuators;
9 
10 namespace UnityStandardAssets.Vehicles.Car
11 {
12     public class CarAgent : Agent
13     {
14         public float spawnRadiusX = 2f;
15         public float spawnRadiusZ = 2f;
16         public float envRadiusX = 3f;
17         public float envRadiusZ = 10f;
18         public float inTargetMultiplier = 1.5f;
19         public GameObject target;
20         public Camera[] carCameras;
21         public int cameraWidth = 200;
22         public int cameraHeight = 200;
23         public bool cameraGrayScale = false;
24         public SensorCompressionType sensorCompressionType = SensorCompressionType.PNG;
25 
26         private CarController carController;
27         EnvironmentParameters defaultParameters;
28         private Rigidbody rb;
29         private int gSteps = 0;
30 
31         private bool inTarget = false;
32 
33         private Vector3 startPosition;
34         private Quaternion startRotation;
35         private Vector3 lastPosition;
36 
37         // Automated parking detection variables
38         public bool findParkingSpot = true;
39         private bool isLookingForSpot;
40         private bool isPositioning;
41         private RayPerceptionSensorComponent3D RayPerceptionSensorComponent;
42         private Vector3 detectedSpotLocation;
43         private float predictedSpotSize = 0f;
```

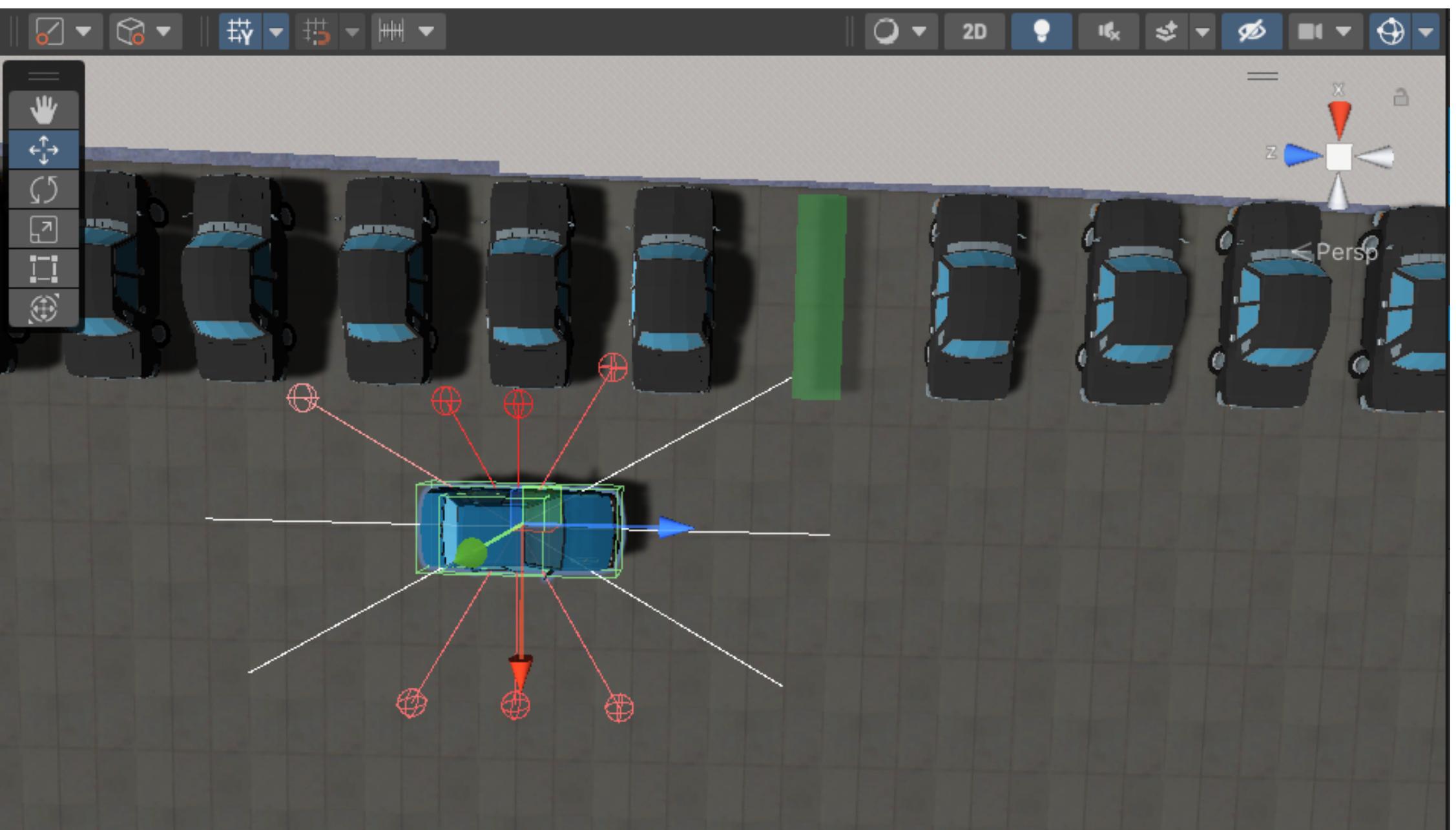
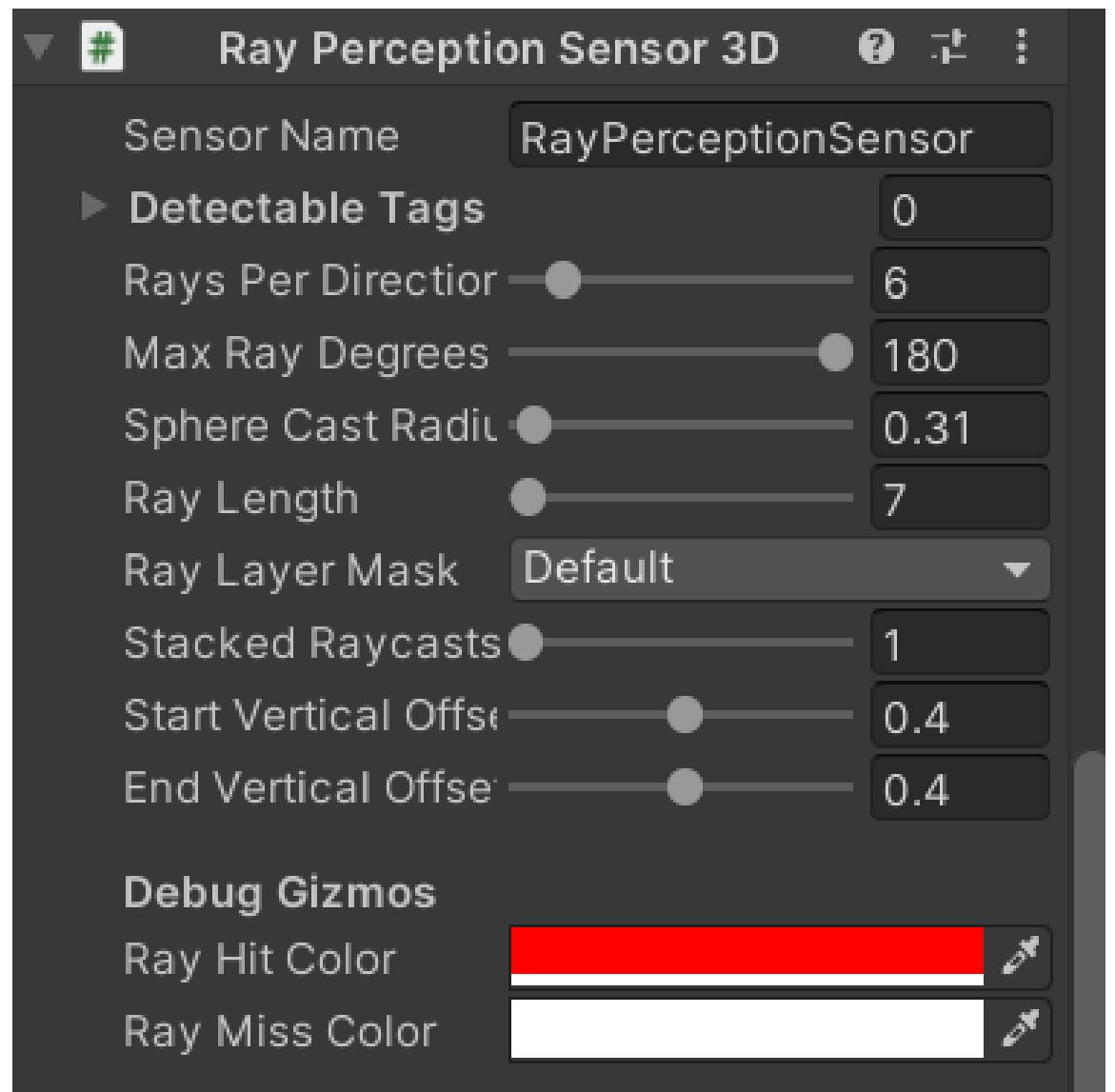
# 場景



The screenshot shows the Unity Asset Store page for the "Japanese Otaku City" asset. The main image is a 3D rendering of a character named Query-Chan, a blonde girl with green eyes wearing a white hoodie. Below the image, there's a video thumbnail showing a city street scene. The title "JAPANESE OTAKU CITY" ASSET PRESENTED BY ZENRIN" is displayed. On the right side, there's a summary section with the title "Japanese Otaku City", the developer "ZENRIN CO., LTD.", a rating of 5 stars (544 reviews), and a "FREE" label. Below that, it says "382 views in the past week". There are "Add to My Assets" and "Love" buttons. A review from "fajuto" is shown, saying "It's just amazing! It's already Unity 3D Models. Unity Assets is greatly upstanding and I am solution you're on you so please last valuable love say was refer everyone...". A "Read more reviews" link is also present. At the bottom, there are links for "License agreement", "Standard Unity Asset Store EULA", "License type", "Extension Asset", "File size", "145.6 MB", and "Latest version", "1.0".

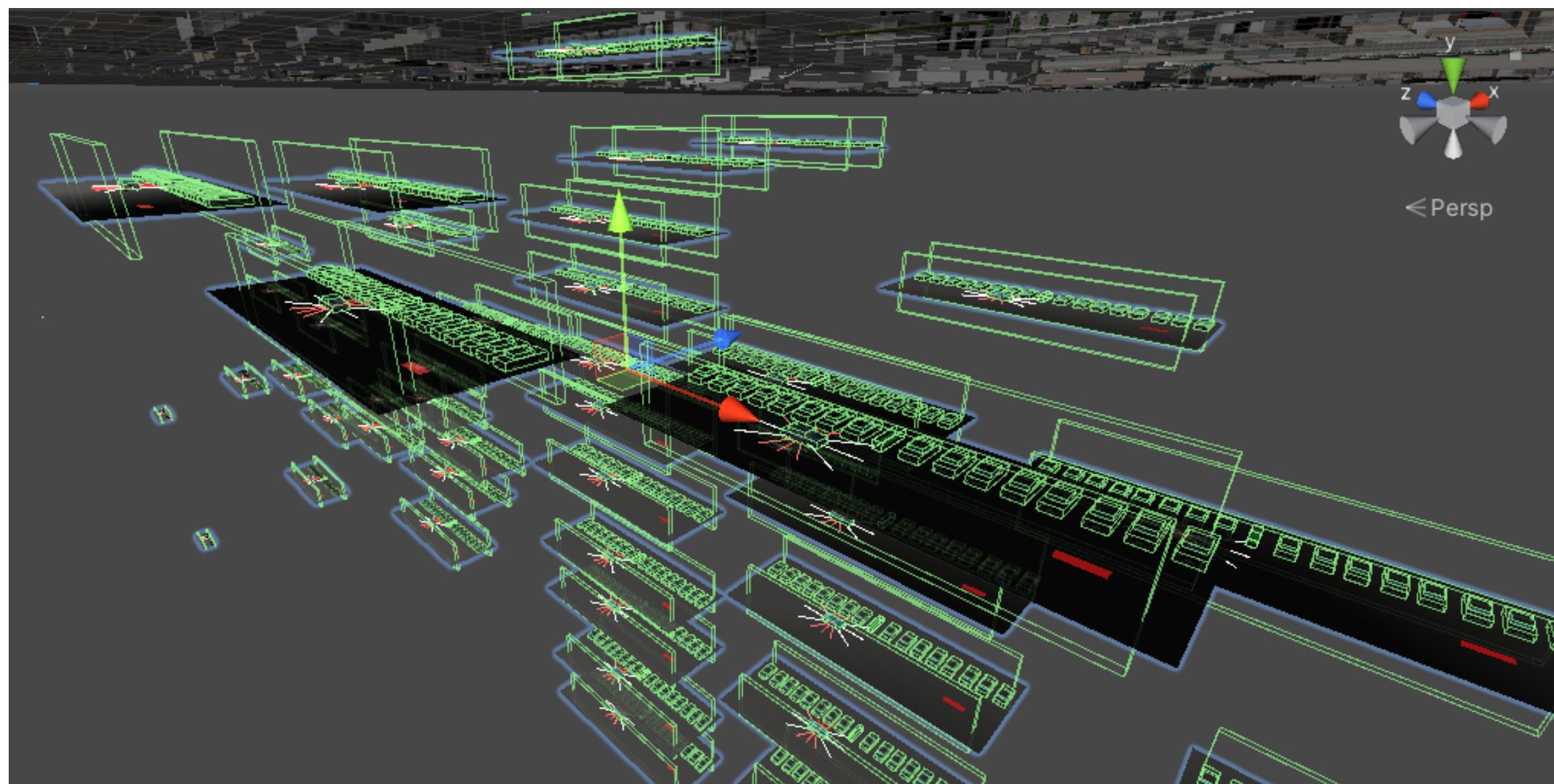


# 傳感器





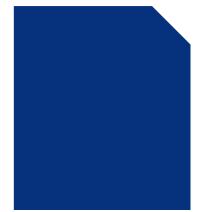
# 訓練陣列 60個環境



# 獎勵算法

# 獎勵機制

增加獎勵:



成功停車(角度、速度)

扣除獎勵:



撞牆、撞車



超出訓練範圍

# 訓練程式

```
2   CarBehaviour:
3     trainer_type: ppo
4     hyperparameters: #超參數設定
5       batch_size: 1024 #每次更新參數時，會有多少步的 (state, action, reward...) 来進行學習
6       buffer_size: 5120 #程式要先集滿5120步數之後才會開啟一輪訓練，再將其分成(5120/1024)個batch，然後每個 batch 進行一次訓練並更新參數，整個過程重複epoch次
7       learning_rate: 0.00035 #梯度下降的學習率
8       beta: 0.0025 #增加熵值來讓動作策略更加隨機，beta越大，越鼓勵探索
9       epsilon: 0.3 #策略更新的速度，越大的話可以使訓練更加穩定，但收斂速度會變更慢
10      lambd: 0.95 #決定了我們的算法多大程度上依賴預估與實際的獎勵值
11      num_epoch: 5 #每次進行梯度下降更新的次數
12      learning_rate_schedule: linear #對學習率使用線性遞減，以達到更快的收斂速度
13
14  network_settings:
15    normalize: true #使用特徵縮放，對於複雜的連續動作才會有幫助
16    hidden_units: 264 #隱藏層的神經元數，數值的大小決定了神經網路對遊戲狀態的表達能力。越多的隱藏層對更大的觀察空間有著較好的表達能力
17    num_layers: 3 #訓練網路的層數。數字越大，模型越深層。但可能會帶來梯度消失與訓練速度慢的問題
18
19  reward_signals:
20    extrinsic: #外在獎勵(環境)
21      gamma: 0.95 #獎勵衰減度
22      strength: 0.99 #獎勵強度
23    gail: #生成對抗模仿學習
24      strength: 0.3 #獎勵強度
25      demo_path: D:/Unity Project/teat1230/Demos/parkteat2.demo #之前訓練出結果的策略文件
26      use_actions: false #不必包含動作，只需策略
27
28  behavioral_cloning: #監督學習，用於限制探索時間，並幫助其探索時找出正確方向。可節省時間與性能
29    demo_path: D:/Unity Project/teat1230/Demos/parkteat2.demo
30    steps: 750000
31    strength: 0.4
32
33  keep_checkpoints: 15 #保留多少個訓練時產生的checkpoint
34  checkpoint_interval: 1000000 #每1000000步後，模型會存一個節點。
35  time_horizon: 264 #在264步數之後，開始把收集到的經驗數據放入到經驗池(類似過擬合、欠擬合)
36  max_steps: 50000000 #最多只會訓練50000000以下步數
37  summary_freq: 100000 #每100000步數之後，開始紀錄我們的訓練統計數據
38  threaded: true #開啟python內建的多線程功能來實現一邊訓練一邊存儲模型
```



**<https://www.youtube.com/watch?v=upyeNufLjs8>**

```
2   CarBehaviour:
3     trainer_type: ppo
4     hyperparameters: #超參數設定
5       batch_size: 1024 #每次更新參數時，會有多少步的 (state, action, reward...) 来進行學習
6       buffer_size: 5120 #程式要先集滿5120步數之後才會開啟一輪訓練，再將其分成(5120/1024)個batch，然後每個 batch 進行一次訓練並更新參數，整個過程重複epoch次
7       learning_rate: 0.00035 #梯度下降的學習率
8       beta: 0.0025 #增加熵值來讓動作策略更加隨機，beta越大，越鼓勵探索
9       epsilon: 0.3 #策略更新的速度，越大的話可以使訓練更加穩定，但收斂速度會變更慢
10      lambd: 0.95 #決定了我們的算法多大程度上依賴預估與實際的獎勵值
11      num_epoch: 5 #每次進行梯度下降更新的次數
12      learning_rate_schedule: linear #對學習率使用線性遞減，以達到更快的收斂速度
13
14  network_settings:
15    normalize: true #使用特徵縮放，對於複雜的連續動作才會有幫助
16    hidden_units: 264 #隱藏層的神經元數，數值的大小決定了神經網路對遊戲狀態的表達能力。越多的隱藏層對更大的觀察空間有著較好的表達能力
17    num_layers: 3 #訓練網路的層數。數字越大，模型越深層。但可能會帶來梯度消失與訓練速度慢的問題
18
19  reward_signals:
20    extrinsic: #外在獎勵(環境)
21      gamma: 0.95 #獎勵衰減度
22      strength: 0.99 #獎勵強度
23    gail: #生成對抗模仿學習
24      strength: 0.3 #獎勵強度
25      demo_path: D:/Unity Project/teat1230/Demos/parkteat2.demo #之前訓練出結果的策略文件
26      use_actions: false #不必包含動作，只需策略
27
28  behavioral_cloning: #監督學習，用於限制探索時間，並幫助其探索時找出正確方向。可節省時間與性能
29    demo_path: D:/Unity Project/teat1230/Demos/parkteat2.demo
30    steps: 750000
31    strength: 0.4
32
33  keep_checkpoints: 15 #保留多少個訓練時產生的checkpoint
34  checkpoint_interval: 1000000 #每1000000步後，模型會存一個節點。
35  time_horizon: 264 #在264步數之後，開始把收集到的經驗數據放入到經驗池(類似過擬合、欠擬合)
36  max_steps: 50000000 #最多只會訓練50000000以下步數
37  summary_freq: 100000 #每100000步數之後，開始紀錄我們的訓練統計數據
38  threaded: true #開啟python內建的多線程功能來實現一邊訓練一邊存儲模型
```

```
[INFO] CarBehaviour. Step: 2800000. Time Elapsed: 3486.212 s. Mean Reward: -109.374. Std of Reward: 170.545. Training.  
[INFO] CarBehaviour. Step: 2900000. Time Elapsed: 3572.751 s. Mean Reward: -93.423. Std of Reward: 175.991. Training.  
[INFO] CarBehaviour. Step: 3000000. Time Elapsed: 3658.638 s. Mean Reward: -108.460. Std of Reward: 179.354. Training.  
[INFO] Exported results\mlagents1\CarBehaviour\CarBehaviour-2999951.onnx  
[INFO] CarBehaviour. Step: 3100000. Time Elapsed: 3759.530 s. Mean Reward: -103.407. Std of Reward: 168.830. Training.  
[INFO] CarBehaviour. Step: 3200000. Time Elapsed: 3855.764 s. Mean Reward: -103.247. Std of Reward: 174.972. Training.  
[INFO] CarBehaviour. Step: 3300000. Time Elapsed: 3950.040 s. Mean Reward: -93.265. Std of Reward: 157.743. Training.  
[INFO] CarBehaviour. Step: 3400000. Time Elapsed: 4040.181 s. Mean Reward: -88.253. Std of Reward: 150.180. Training.  
[INFO] CarBehaviour. Step: 3500000. Time Elapsed: 4137.314 s. Mean Reward: -100.797. Std of Reward: 172.829. Training.  
[INFO] CarBehaviour. Step: 3600000. Time Elapsed: 4232.873 s. Mean Reward: -97.138. Std of Reward: 169.860. Training.  
[INFO] CarBehaviour. Step: 3700000. Time Elapsed: 4321.822 s. Mean Reward: -106.370. Std of Reward: 185.393. Training.  
[INFO] CarBehaviour. Step: 3800000. Time Elapsed: 4414.687 s. Mean Reward: -109.490. Std of Reward: 169.717. Training.  
[INFO] CarBehaviour. Step: 3900000. Time Elapsed: 4504.885 s. Mean Reward: -105.258. Std of Reward: 157.714. Training.  
[INFO] CarBehaviour. Step: 4000000. Time Elapsed: 4598.343 s. Mean Reward: -101.874. Std of Reward: 165.826. Training.  
[INFO] Exported results\mlagents1\CarBehaviour\CarBehaviour-3999938.onnx  
[INFO] CarBehaviour. Step: 4100000. Time Elapsed: 4691.743 s. Mean Reward: -94.384. Std of Reward: 173.884. Training.  
[INFO] CarBehaviour. Step: 4200000. Time Elapsed: 4785.094 s. Mean Reward: -100.277. Std of Reward: 182.768. Training.  
[INFO] CarBehaviour. Step: 4300000. Time Elapsed: 4878.896 s. Mean Reward: -94.410. Std of Reward: 147.766. Training.  
[INFO] CarBehaviour. Step: 4400000. Time Elapsed: 4971.541 s. Mean Reward: -97.918. Std of Reward: 167.456. Training.  
[INFO] CarBehaviour. Step: 4500000. Time Elapsed: 5068.152 s. Mean Reward: -83.488. Std of Reward: 151.603. Training.  
[INFO] CarBehaviour. Step: 4600000. Time Elapsed: 5160.794 s. Mean Reward: -100.727. Std of Reward: 193.059. Training.  
[INFO] CarBehaviour. Step: 4700000. Time Elapsed: 5255.509 s. Mean Reward: -80.636. Std of Reward: 200.708. Training.  
[INFO] CarBehaviour. Step: 4800000. Time Elapsed: 5346.742 s. Mean Reward: -105.449. Std of Reward: 223.087. Training.  
[INFO] CarBehaviour. Step: 4900000. Time Elapsed: 5438.903 s. Mean Reward: -71.209. Std of Reward: 197.954. Training.  
[INFO] CarBehaviour. Step: 5000000. Time Elapsed: 5536.375 s. Mean Reward: -83.424. Std of Reward: 188.374. Training.  
[INFO] Exported results\mlagents1\CarBehaviour\CarBehaviour-4999830.onnx  
[INFO] CarBehaviour. Step: 5100000. Time Elapsed: 5633.252 s. Mean Reward: -99.596. Std of Reward: 192.804. Training.  
[INFO] CarBehaviour. Step: 5200000. Time Elapsed: 5728.199 s. Mean Reward: -84.441. Std of Reward: 205.504. Training.  
[INFO] CarBehaviour. Step: 5300000. Time Elapsed: 5822.570 s. Mean Reward: -87.887. Std of Reward: 178.291. Training.  
[INFO] CarBehaviour. Step: 5400000. Time Elapsed: 5916.263 s. Mean Reward: -67.072. Std of Reward: 171.992. Training.  
[INFO] CarBehaviour. Step: 5500000. Time Elapsed: 6011.229 s. Mean Reward: -75.982. Std of Reward: 159.282. Training.  
[INFO] CarBehaviour. Step: 5600000. Time Elapsed: 6105.291 s. Mean Reward: -94.287. Std of Reward: 178.745. Training.  
[INFO] CarBehaviour. Step: 5700000. Time Elapsed: 6197.936 s. Mean Reward: -107.858. Std of Reward: 185.582. Training.  
[INFO] CarBehaviour. Step: 5800000. Time Elapsed: 6293.844 s. Mean Reward: -96.976. Std of Reward: 191.394. Training.  
[INFO] CarBehaviour. Step: 5900000. Time Elapsed: 6388.978 s. Mean Reward: -87.319. Std of Reward: 179.422. Training.  
[INFO] CarBehaviour. Step: 6000000. Time Elapsed: 6482.283 s. Mean Reward: -80.371. Std of Reward: 176.388. Training.  
[INFO] Exported results\mlagents1\CarBehaviour\CarBehaviour-5999909.onnx  
[INFO] CarBehaviour. Step: 6100000. Time Elapsed: 6572.174 s. Mean Reward: -92.133. Std of Reward: 162.141. Training.  
[INFO] CarBehaviour. Step: 6200000. Time Elapsed: 6671.574 s. Mean Reward: -82.007. Std of Reward: 153.439. Training.  
[INFO] CarBehaviour. Step: 6300000. Time Elapsed: 6766.448 s. Mean Reward: -74.190. Std of Reward: 204.158. Training.
```



# Mean reward

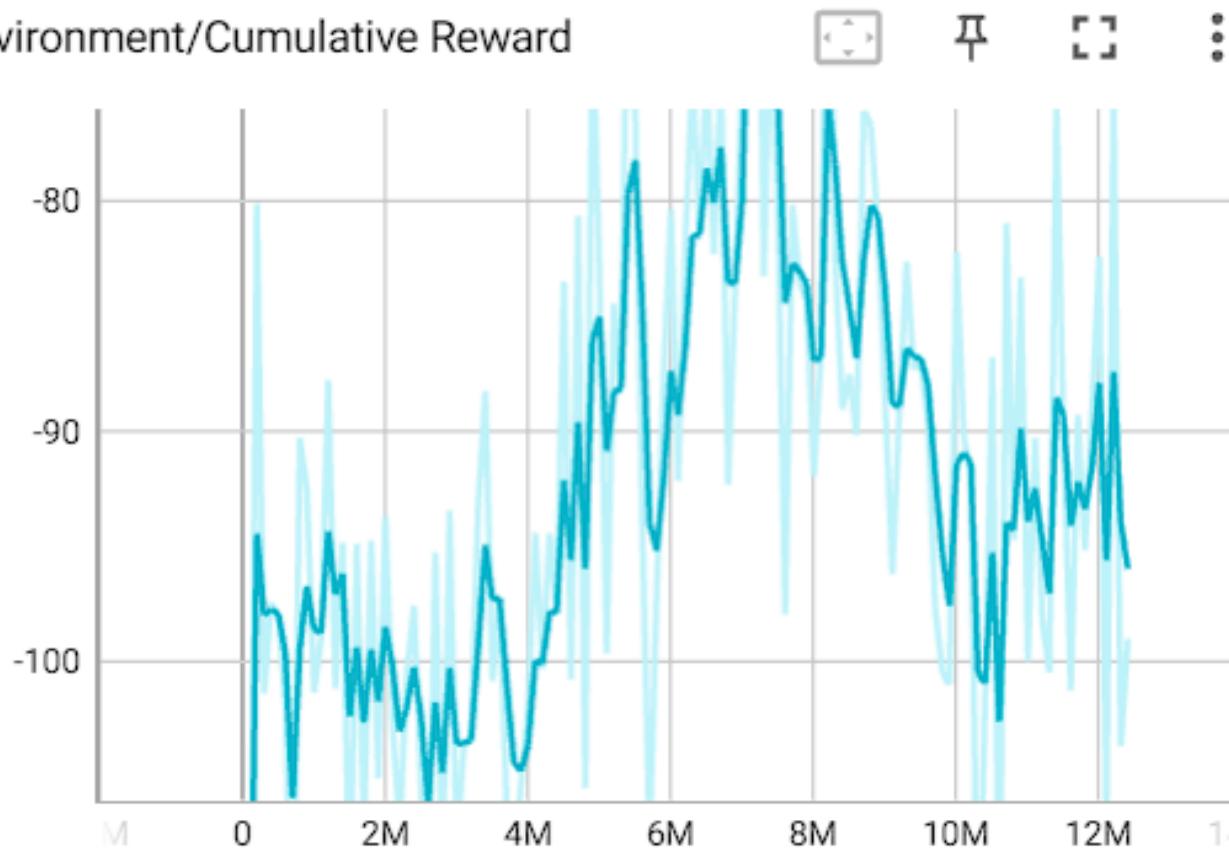
在每10萬步裡, Agents 拿到的Reward

# Std reward

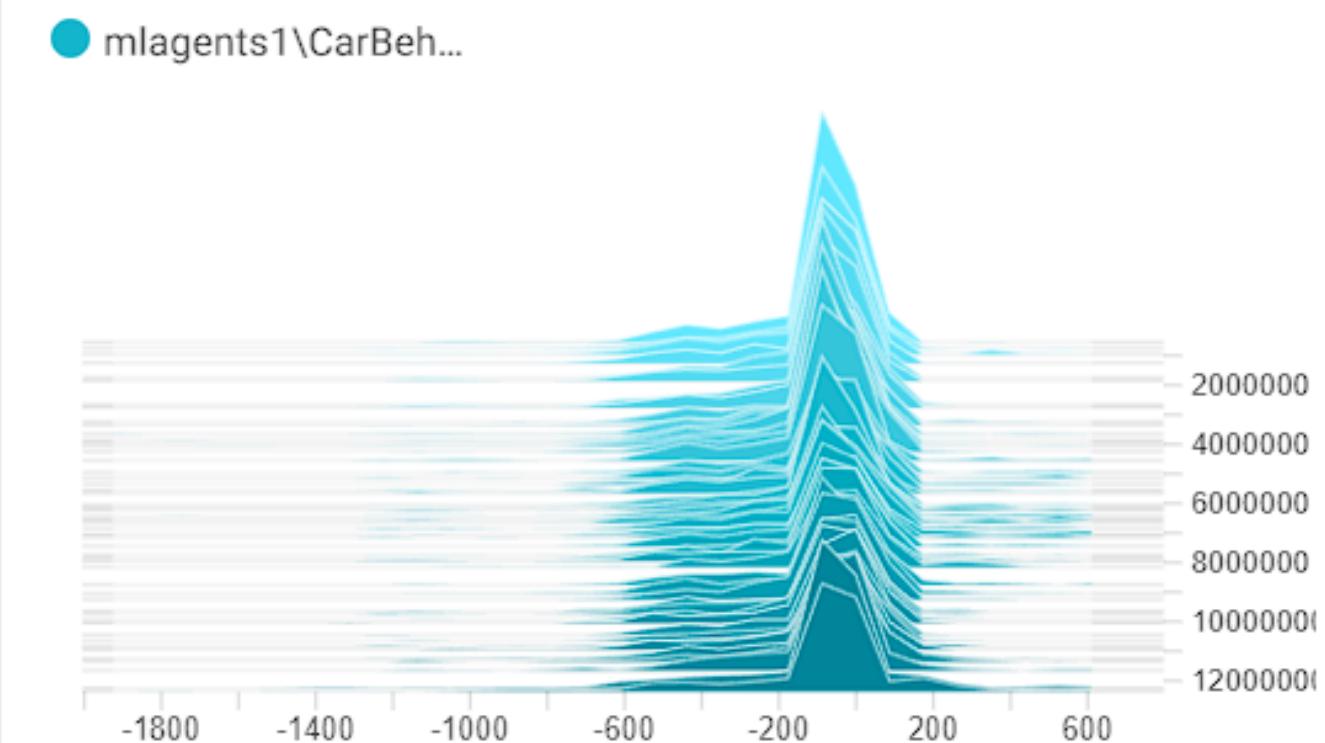
衡量平均獎勵分佈的量度。較大的值表示收到的獎勵變化很大，而較小的值則相反。

# Tensorboard

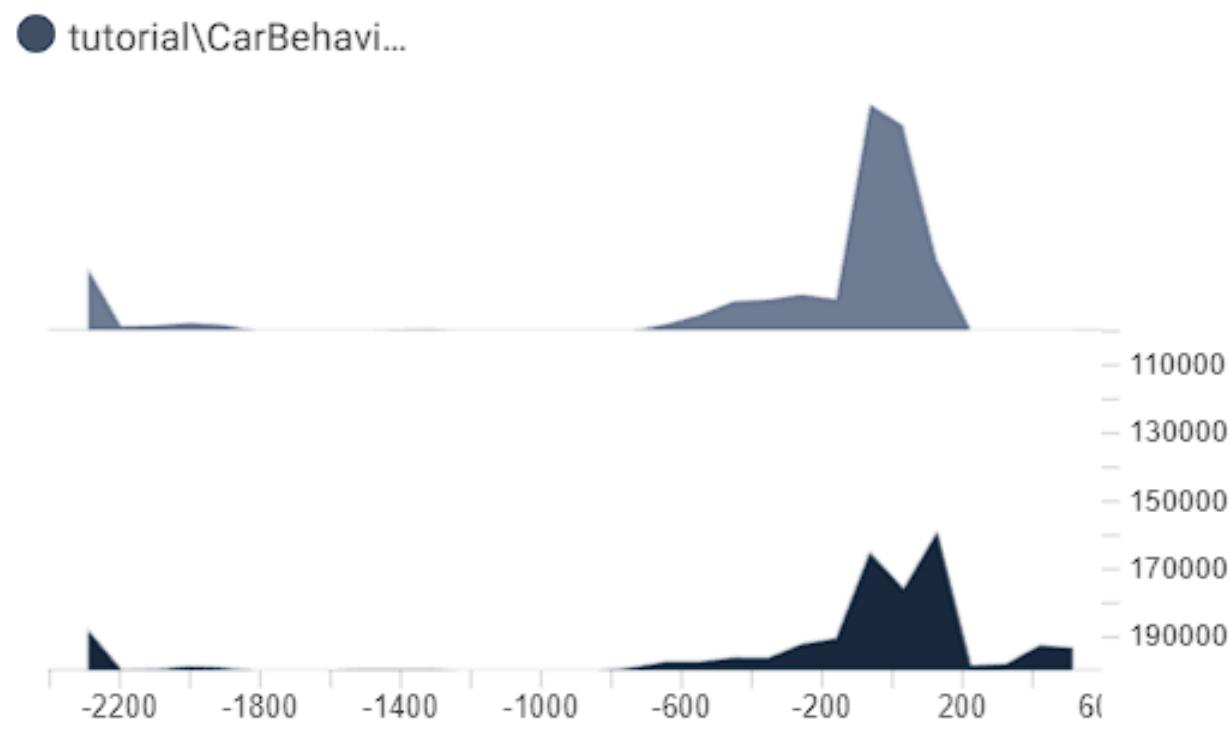
Environment/Cumulative Reward



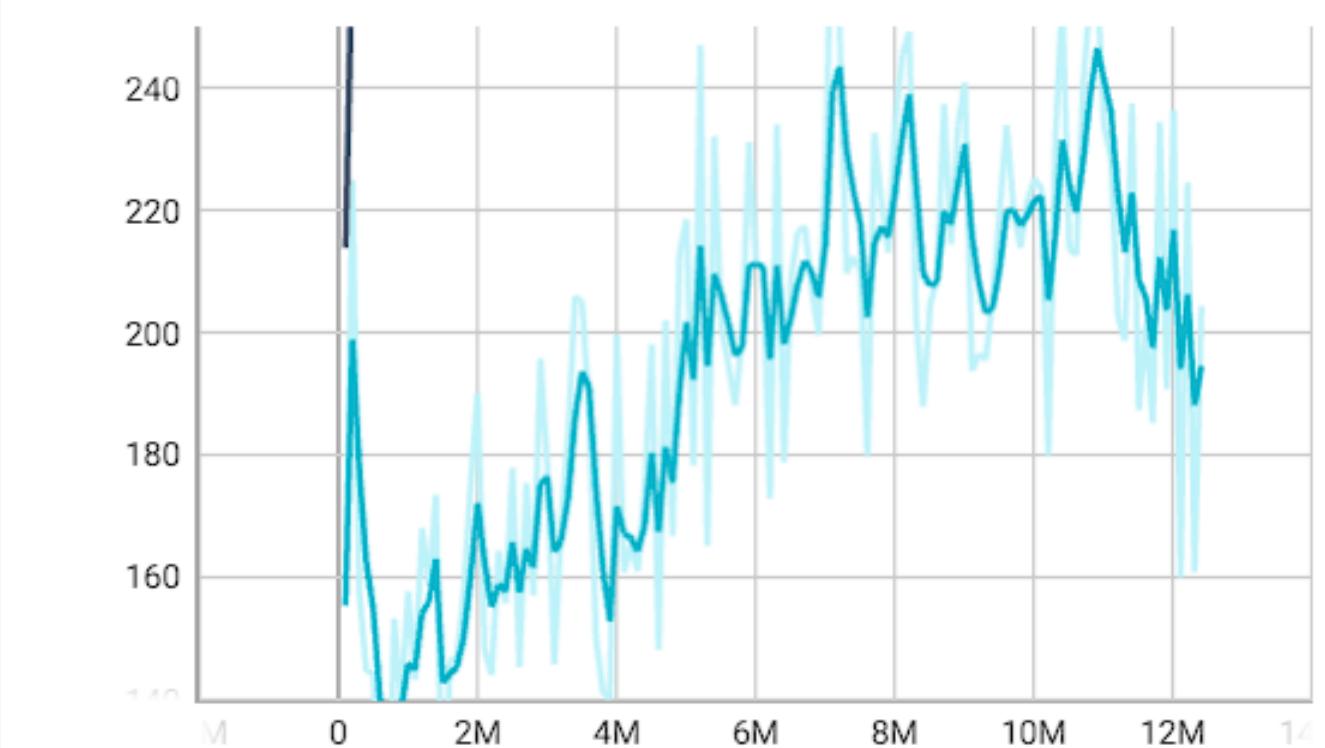
Environment/Cumulative Reward\_hist



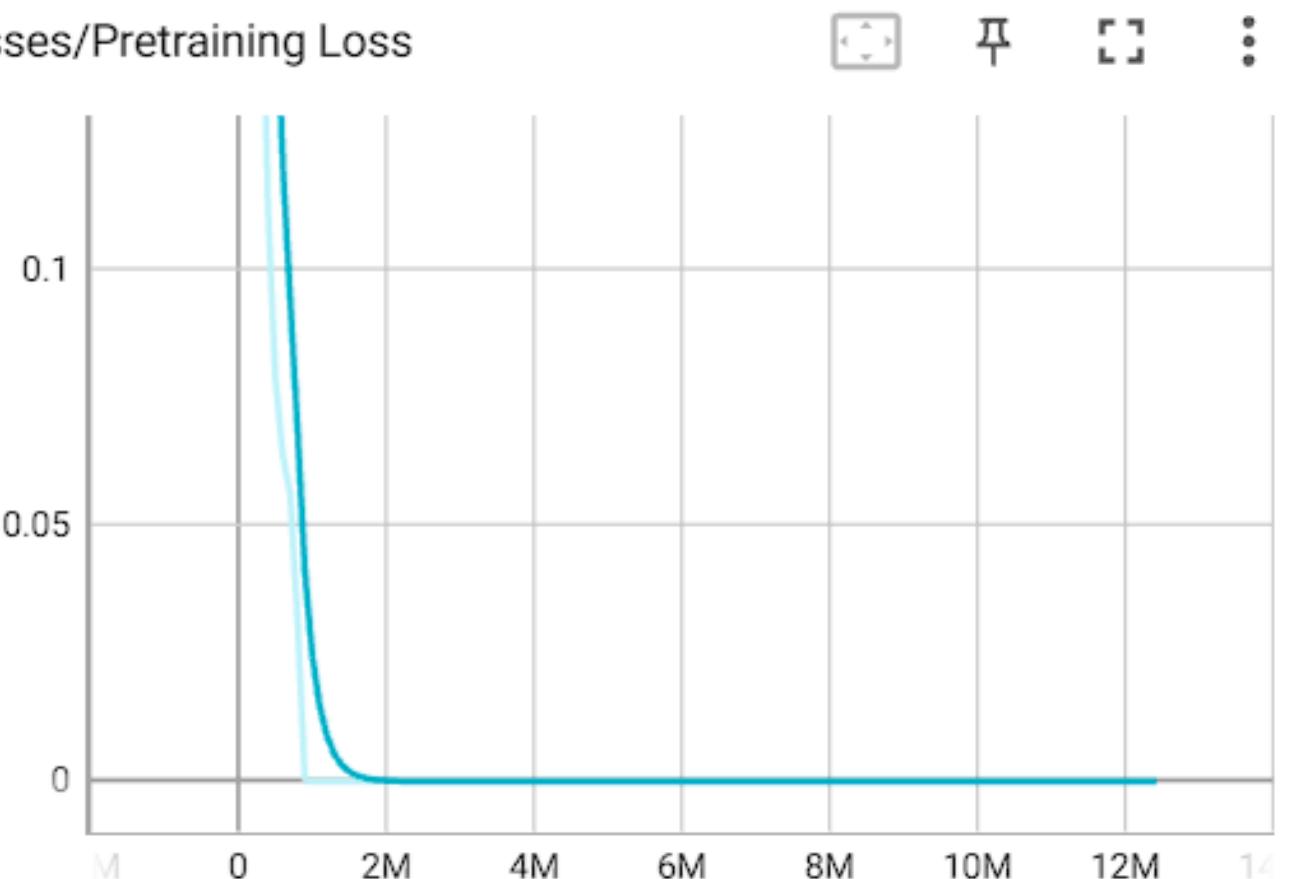
Environment/Cumulative Reward\_hist



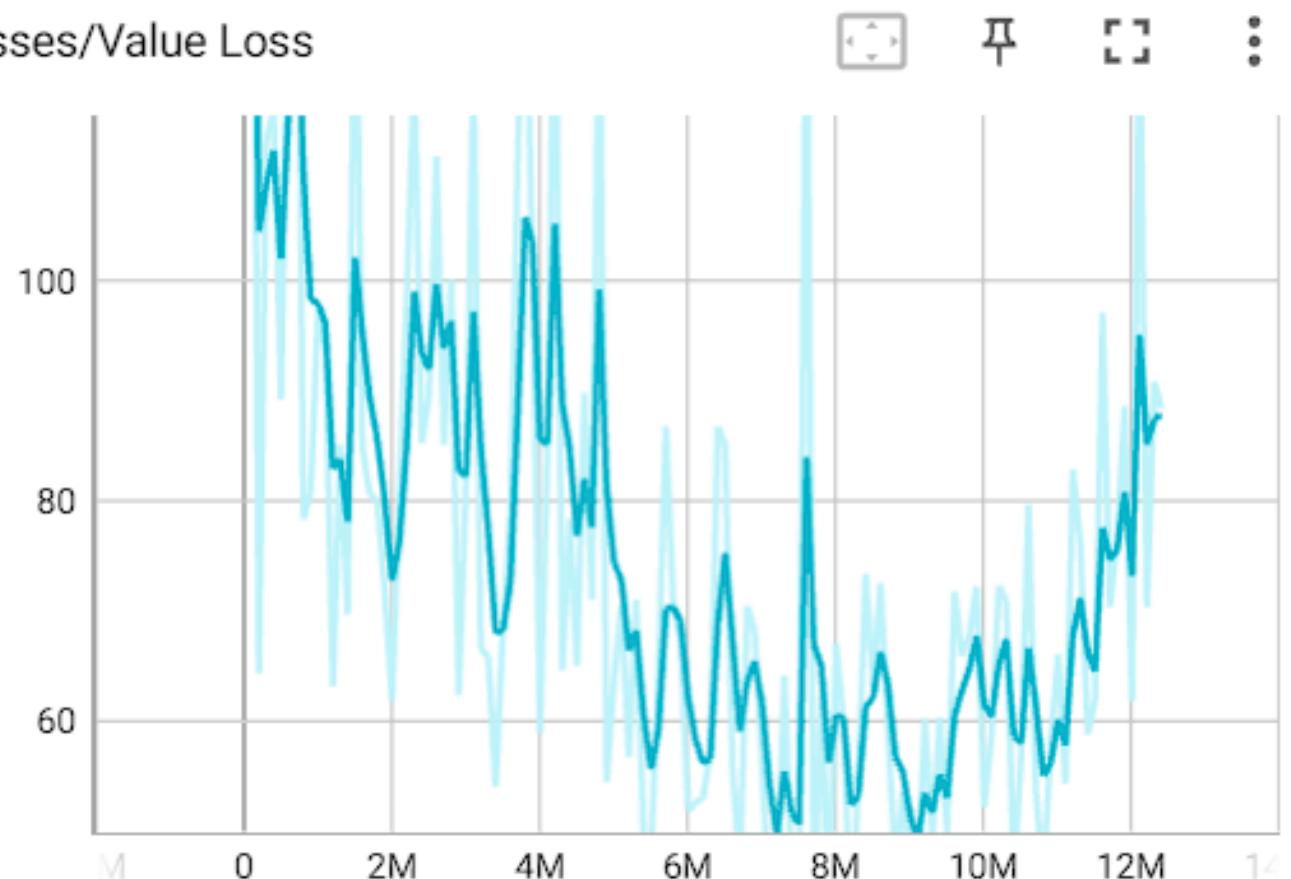
Environment/Episode Length



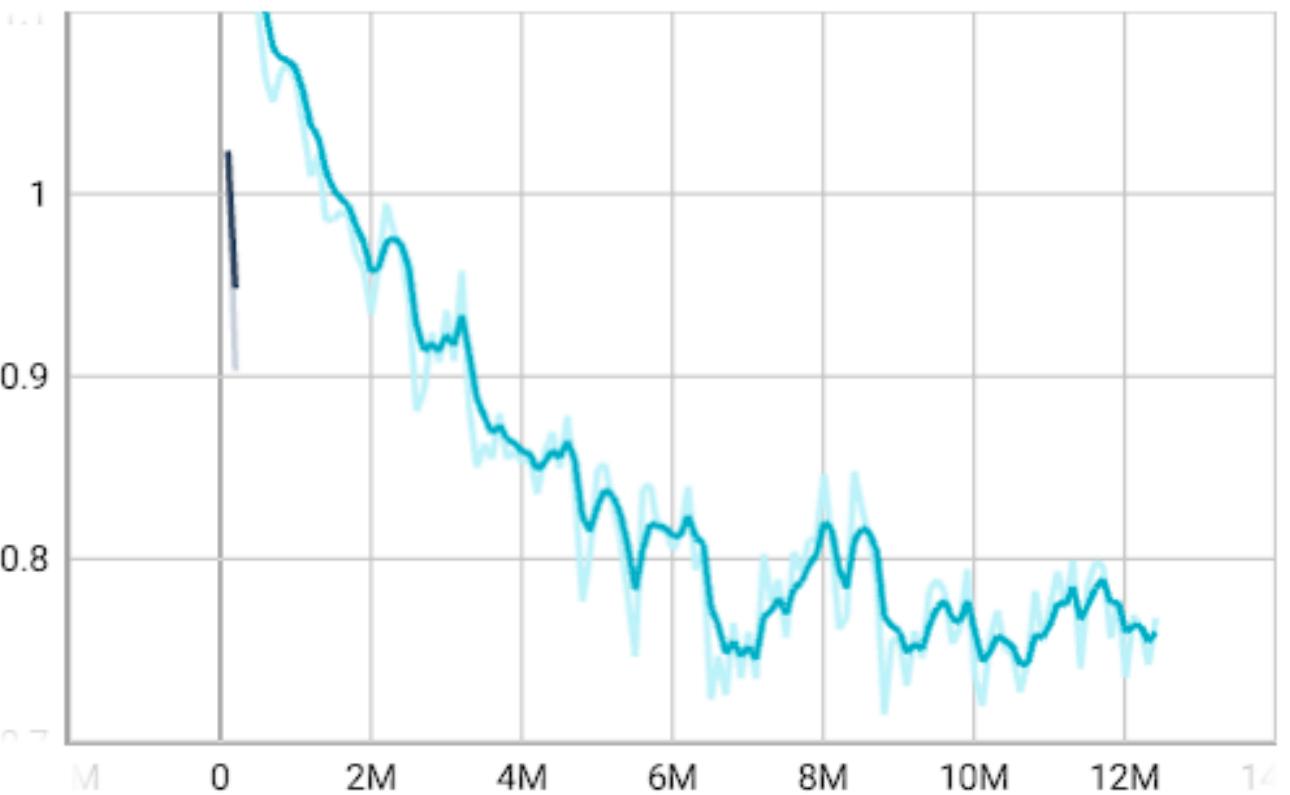
Losses/Pretraining Loss



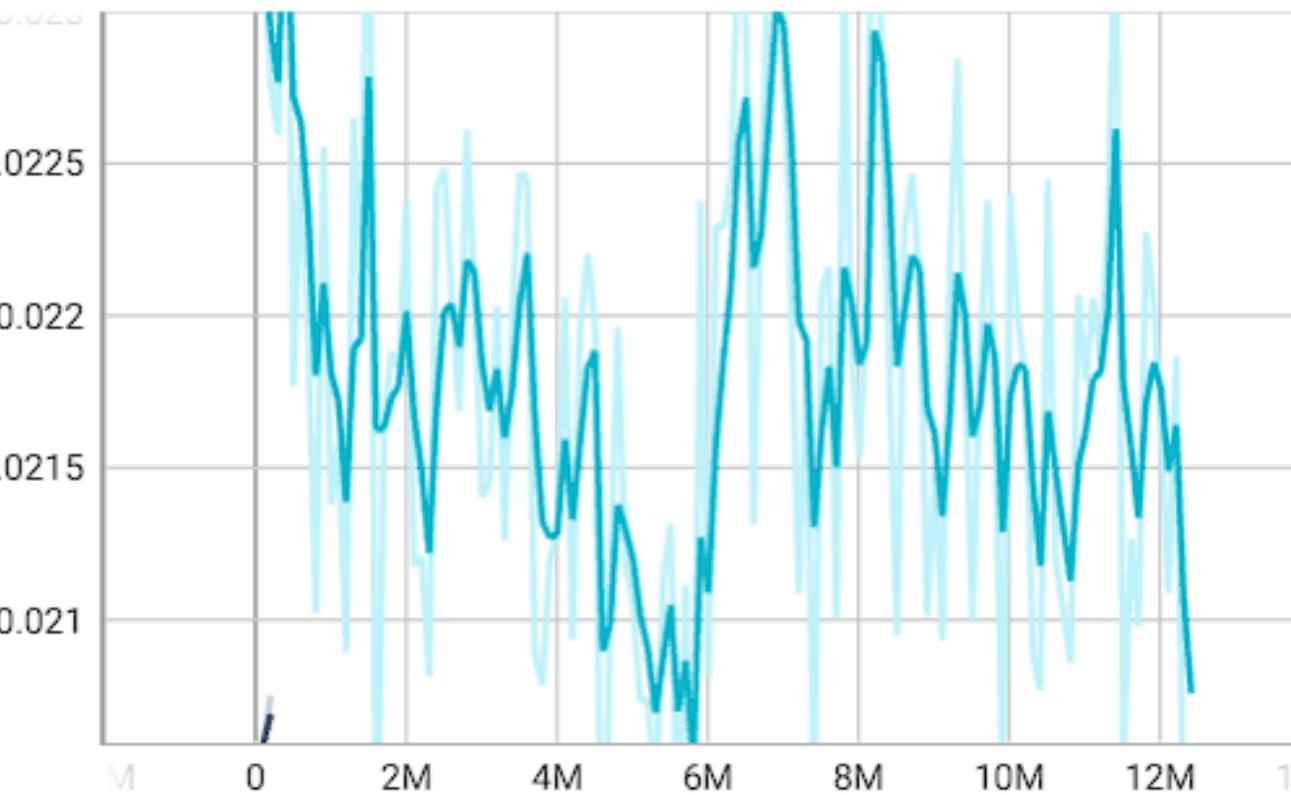
Losses/Value Loss



Losses/GAIL Loss



Losses/Policy Loss



# 成果展示

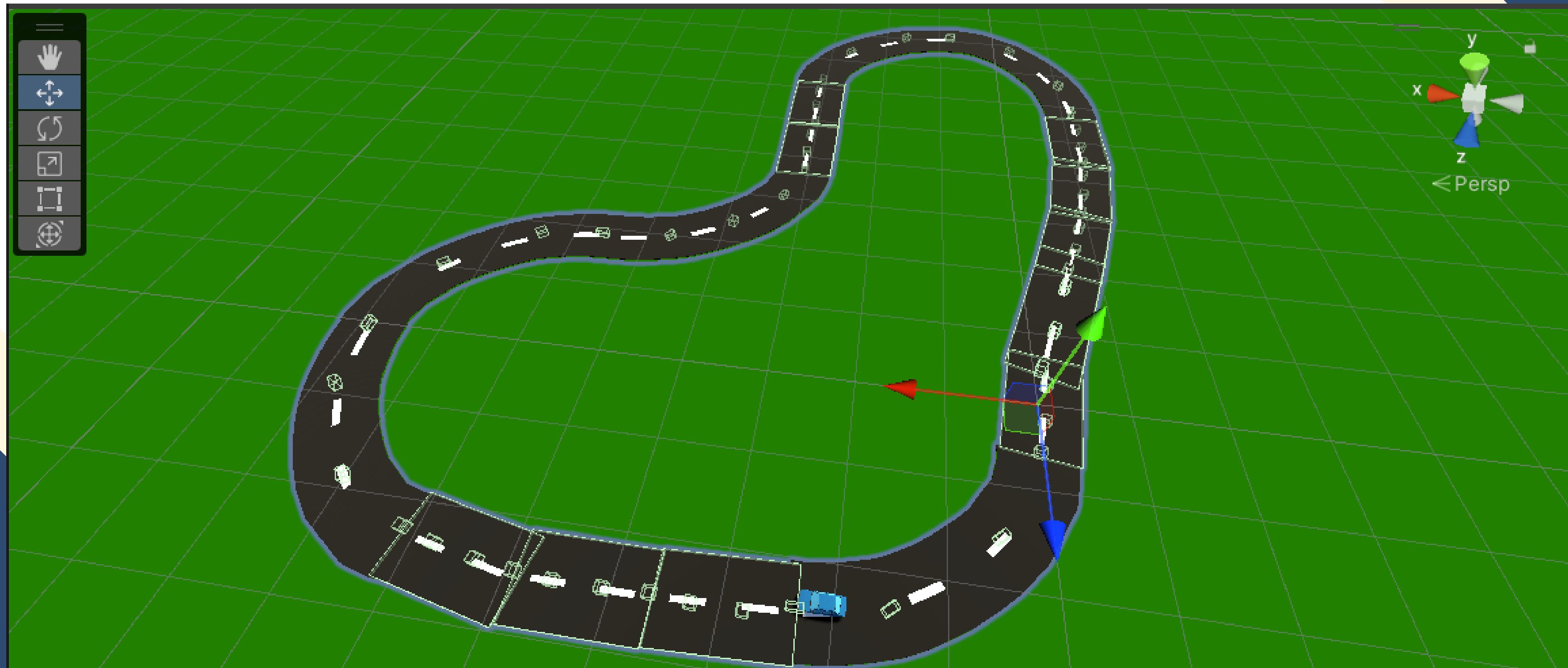
**<https://www.youtube.com/watch?v=KmknP3KVYIs>**

# 貢獻與後續

# 貢獻

1. 更改環境，更改與優化參數
2. Tensorboard 分析
3. 優化car spawner
4. 自己的模型
5. 自走車半成品(無法轉彎)

# 自走車原理



# 如何改善

訓練跑滿5千萬步  
更好的獎勵機制

**謝謝各位**