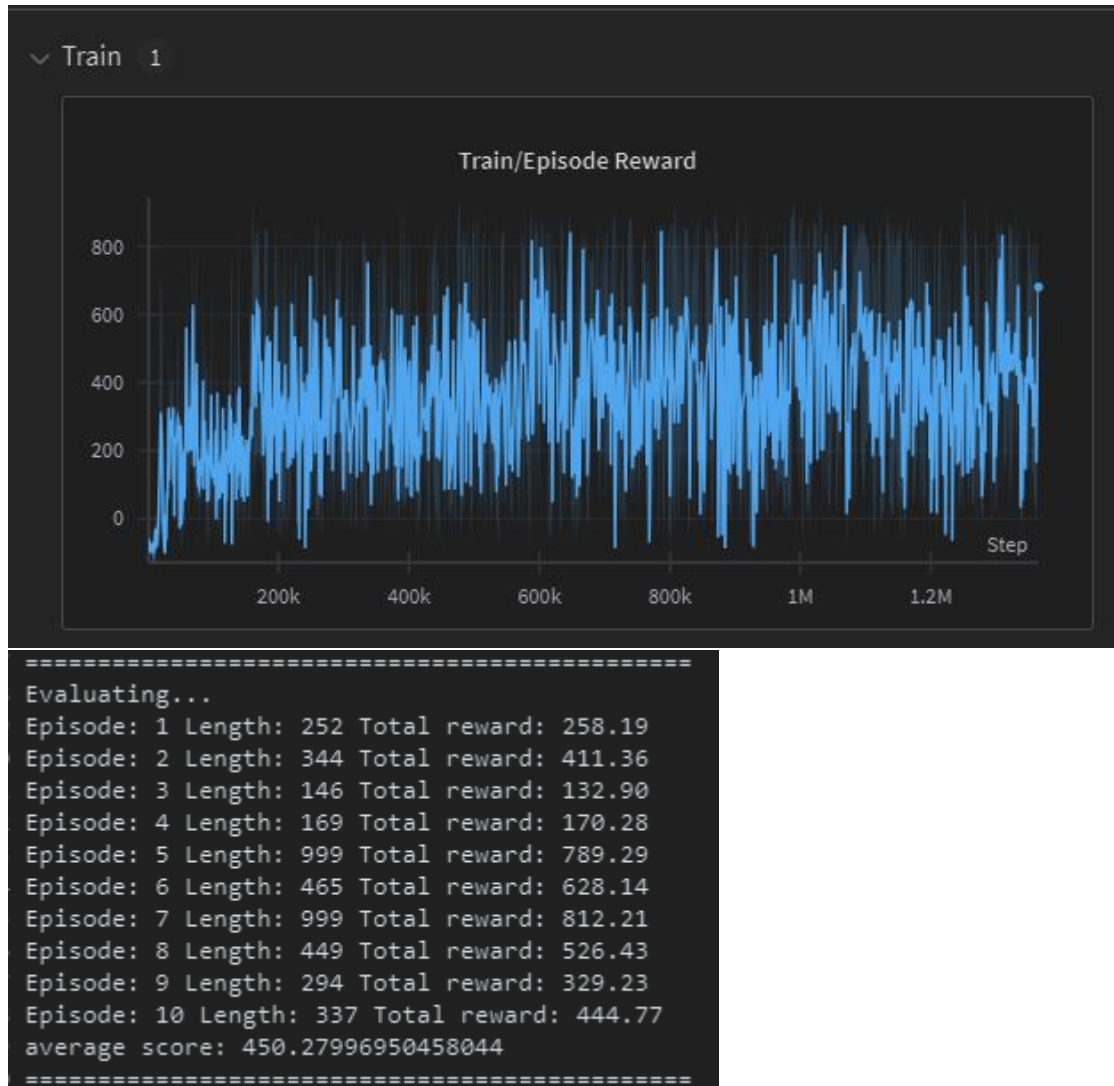


## RL Topic HW3

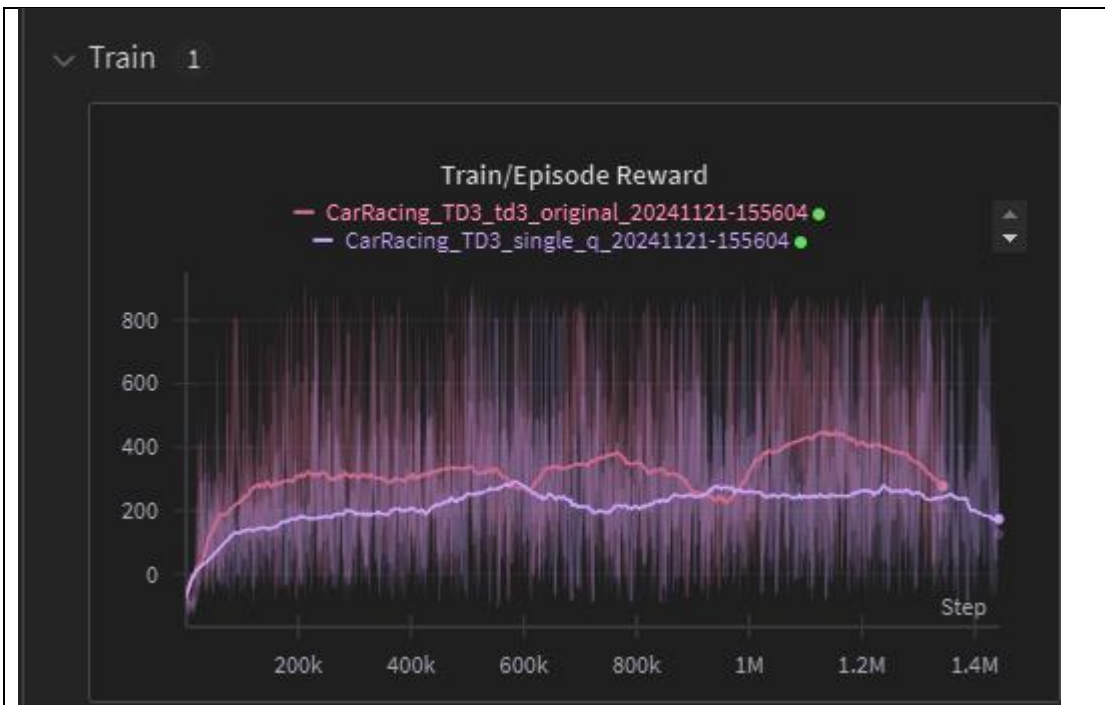
### 1. Training curve and testing result of TD3:



### 2. Questions:

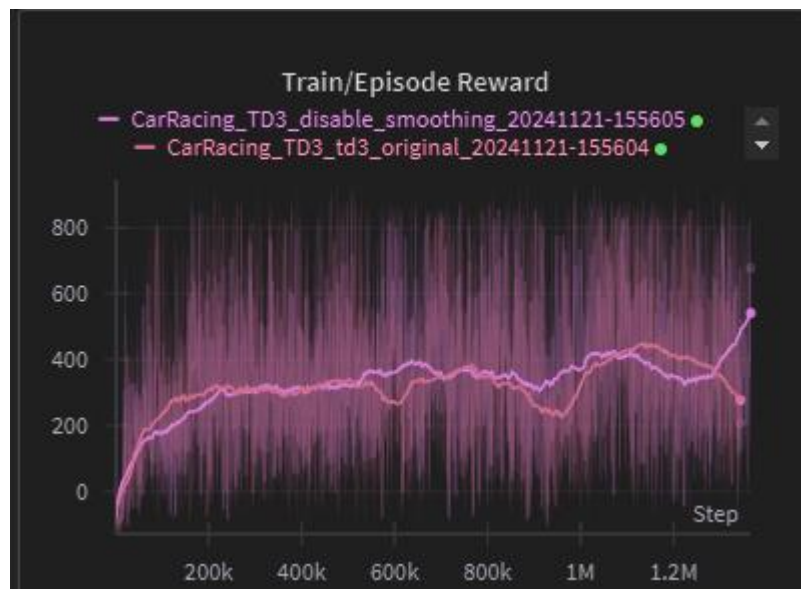
#### a. Single vs Double Q network

Twin Q-Networks in TD3 improve performance by reducing overestimation bias through the use of the minimum Q-value from two networks, leading to more accurate value estimates and stable learning. In contrast, a Single Q-Network is simpler but prone to overestimation, which can result in suboptimal policies and less stable training. Twin Q-networks trade increased computational cost for better reliability and performance.



b. Disable target policy smoothing

Enabling target policy smoothing in TD3 improves stability by adding noise to target actions, making the Q-network less sensitive to small action errors and reducing overfitting to deterministic policies. Disabling it can lead to less robust policies, as the Q-network may overfit to specific actions, making training less stable in noisy environments.



c. Disable delayed update

Delayed update steps in TD3 improve stability by updating the policy and target networks less frequently than the Q-networks. This prevents the policy from being overly influenced by noisy or unstable Q-value estimates. Without delayed updates, the policy might overfit

to fluctuating Q-values, leading to instability and poorer performance. Delayed updates trade slower updates for more robust and reliable learning.

