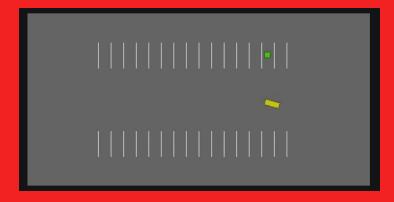
Deep Reinforcement Learning Implementation on a Parking Environment

- Team Members:
 - Aditya Aspat
 - Atharva Jamsandekar







Problem Statement & Environment

- Objective: To park the Agent in the highlighted parking spot in the parallel parking orientation.
 - Continuous action spaces with multiple actions like steering and acceleration pose challenges due to high dimensionality, requiring precise control, complex policy representation, and effective exploration strategies in reinforcement learning tasks.

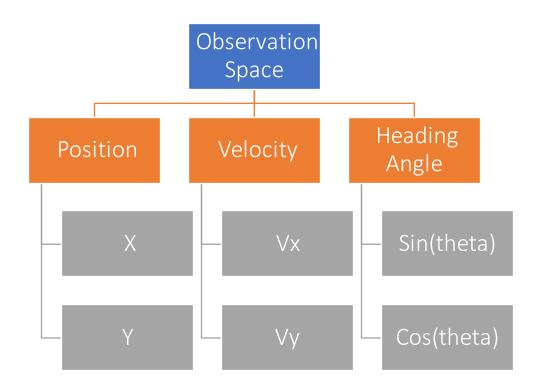
Motivation:

- Action Space is identical to the control space of real-world Vehicles.
- The domain is a stepping stone to understand the research in the continuous action space domain.

Rewards:

- The Lp norm is calculated between the achieved state and the desired state after multiplying with a reward weights vector.
- Collision Reward: -5
- Success Reward: 0.12

$$\ell_p = \left(\sum_{i=1}^N \left|x_i
ight|^p
ight)^{1/p}, ext{for } p \geq 1$$





Deep Deterministic Policy Gradient (DDPG)

Features:

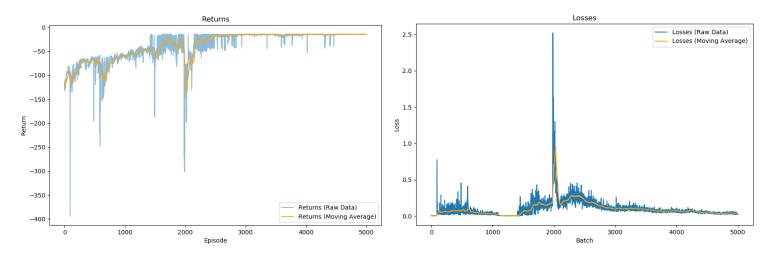
- Actor-Critic with Target Network
- Continuous Action & Observation Space
- Maximizes the expected cumulative long-term reward
- Off Policy Algorithm
- Replay Buffer
- Soft Update of Target Networks

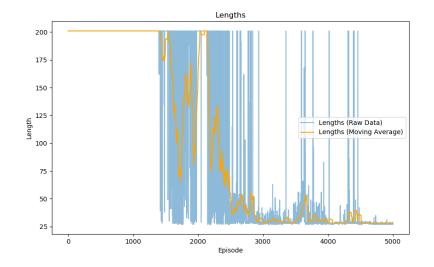
$$\begin{split} J_{Q} &= \frac{1}{N} \sum_{i=1}^{N} (r_{i} + \gamma(1 - d)Q_{targ}(s_{i}', \mu_{targ}(s_{i}')) - Q(s_{i}, \mu(s_{i}))^{2} \\ J_{\mu} &= \frac{1}{N} \sum_{i=1}^{N} Q(s_{i}, \mu(s_{i})) \end{split} \qquad \begin{aligned} \theta^{\mu}_{targ} &\leftarrow \tau \theta^{\mu}_{targ} + (1 - \tau)\theta^{\mu} \\ \theta^{Q}_{targ} &\leftarrow \tau \theta^{Q}_{targ} + (1 - \tau)\theta^{Q} \end{aligned}$$

Hyperparameters:

- OU Noise (mu=0, theta=0.15, sigma=0.2) $dx_t = \theta(\mu x_t)dt + \sigma dW_t$
- Gamma=0.99
- Tau=0.005

Results:





Proximal Policy Optimization (PPO)

Features:

- Actor-Critic Structure
- Continuous Action & Observation Space
- On Policy Algorithm
- Replay Buffer
- Clipped Objective function.

$$\mathcal{L}^{\textit{CLIP}}_{ heta_k}(heta) = \mathop{\mathbb{E}}_{ au \sim \pi_k} \left[\sum_{t=0}^T \left[\min(r_t(heta) \hat{A}^{\pi_k}_t, \operatorname{clip}\left(r_t(heta), 1 - \epsilon, 1 + \epsilon
ight) \hat{A}^{\pi_k}_t
ight)
ight]
ight]$$

Hyperparameters:

• K_epochs = 40 (update policy for K epochs)

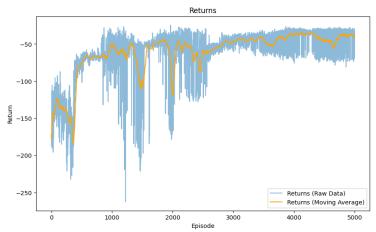
eps_clip = 0.2 (clip parameter for PPO)

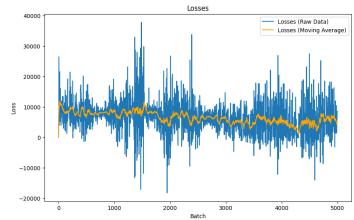
gamma = 0.99 (discount factor)

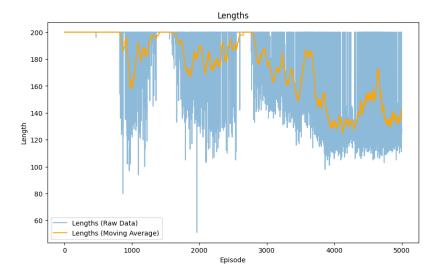
lr_actor = 0.0003 (learning rate for actor network)

• Ir_critic = 0.001 (learning rate for critic network)

Results:







object File]

Features:

- 1 Actor 2 Critic Network Structure
- Continuous Action & Observation Space
- Maximize the entropy of the policy
- Off Policy Algorithm
- Replay Buffer
- Soft Update of Target Networks

$$J(\pi) = \sum_{t=0}^{T} \mathbb{E}_{(\mathbf{s}_{t}, \mathbf{a}_{t}) \sim \rho_{\pi}} \left[r(\mathbf{s}_{t}, \mathbf{a}_{t}) + \alpha \mathcal{H}(\pi(\cdot | \mathbf{s}_{t})) \right].$$

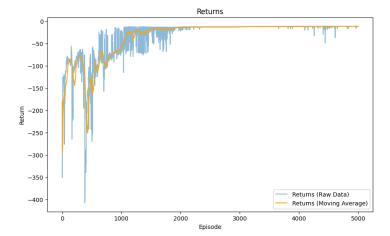
$$\hat{\nabla}_{\psi} J_{V}(\psi) = \nabla_{\psi} V_{\psi}(\mathbf{s}_{t}) \left(V_{\psi}(\mathbf{s}_{t}) - Q_{\theta}(\mathbf{s}_{t}, \mathbf{a}_{t}) + \log \pi_{\phi}(\mathbf{a}_{t} | \mathbf{s}_{t}) \right)$$

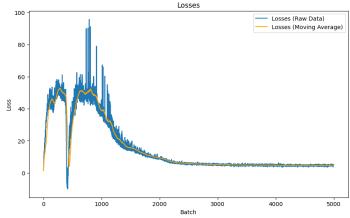
$$J_Q(\theta) = \mathbb{E}_{(\mathbf{s}_t, \mathbf{a}_t) \sim \mathcal{D}} \left[\frac{1}{2} \left(Q_{\theta}(\mathbf{s}_t, \mathbf{a}_t) - \hat{Q}(\mathbf{s}_t, \mathbf{a}_t) \right)^2 \right]$$

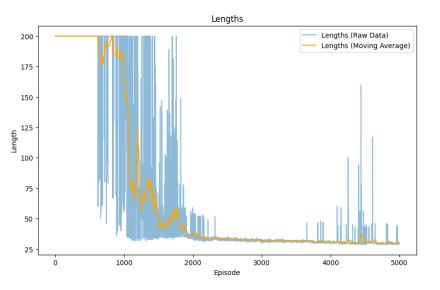
• Hyperparameters:

- alpha = 0.2
- gamma = 0.99
- tau = 0.005
- Learning rate= 0.0003
- batch_size = 64

Results:







object File]

- DDPG performs the best amongst the implemented algorithms.
- PPO has decent performance that gives near-optimal results.
- Feature Association was implemented but it gave suboptimal results.
- SAC converges slower than others but has good performance.
- Future Work:
 - Try to improve the results in Feature Association Model.
 - Apply the implementations to similar domains in our fields of research interests.