

# Learning Safe Cooperative Policies in Autonomous Multi-UAV Navigation

Arshdeep Singh

Indian Institute of Technology Ropar

[2019csm1001@iitrpr.ac.in](mailto:2019csm1001@iitrpr.ac.in)

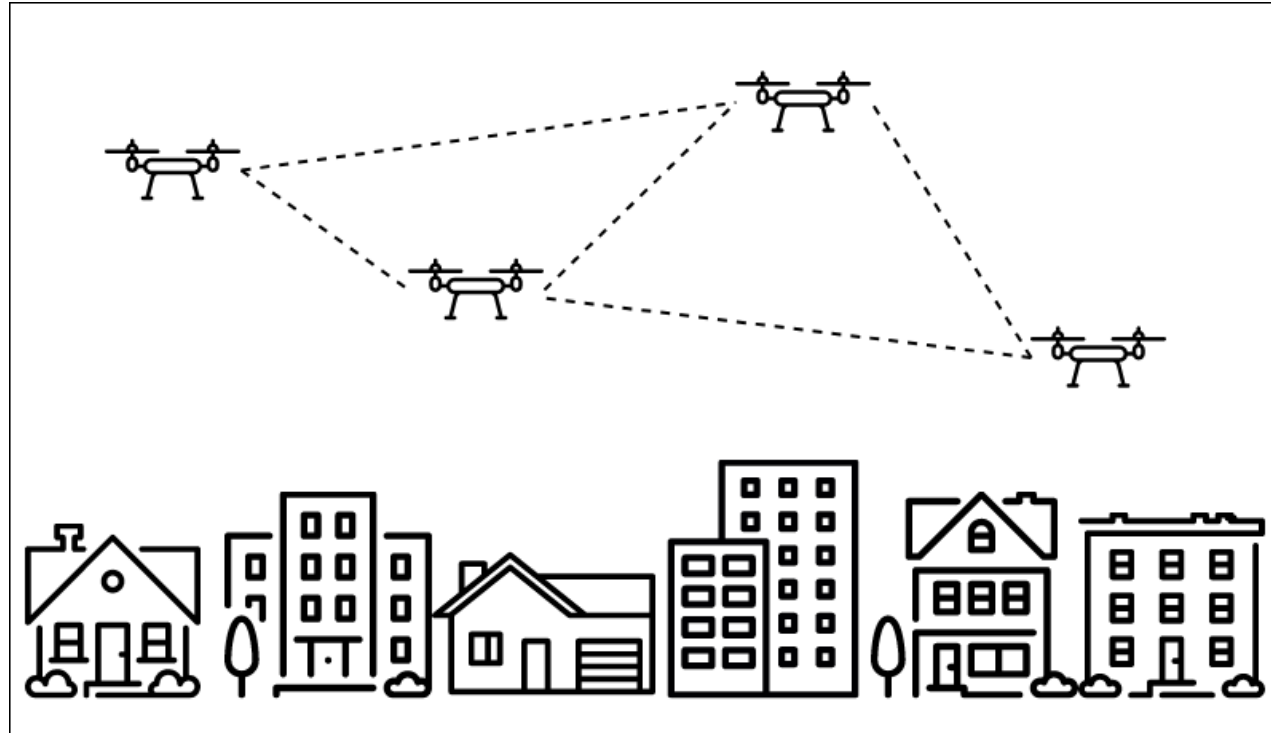
Shashi Shekhar Jha

Indian Institute of Technology Ropar

[shashi@iitrpr.ac.in](mailto:shashi@iitrpr.ac.in)

# Multi-UAV Applications

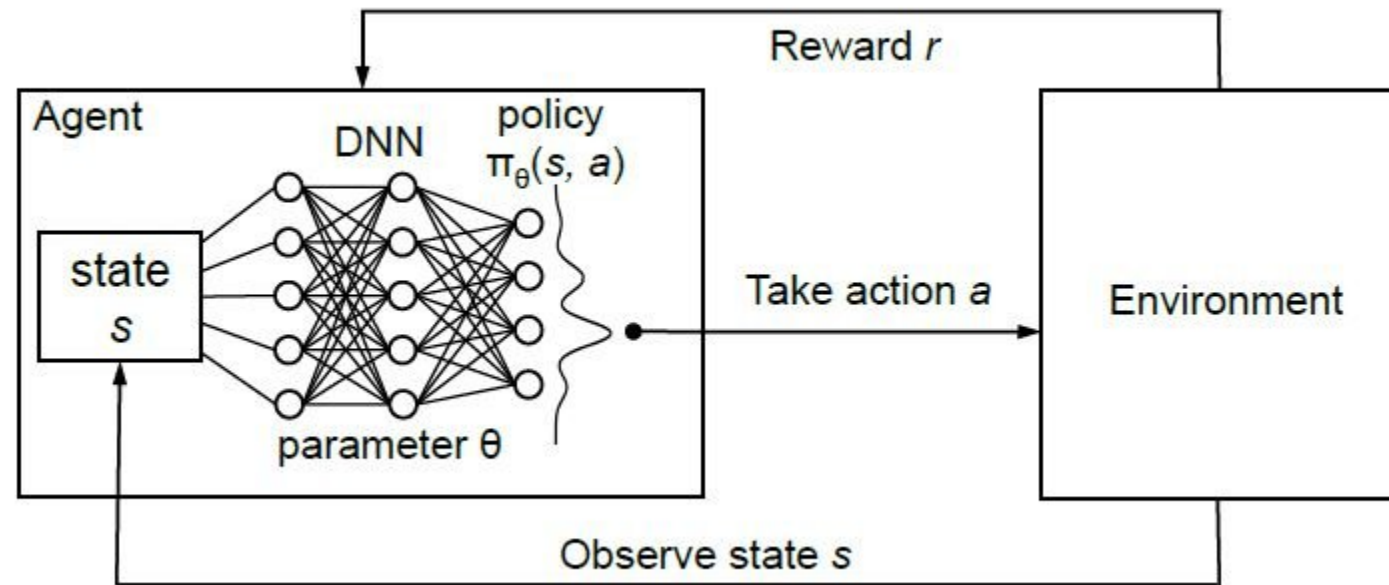
- Application Domains
  - Surveying, Remote Sensing, Search & Rescue, E-Commerce Delivery



# Problem Formulation

- Multi-UAVs learning to autonomously reach their target destination through safe navigation
- Formulated as a Multi Agent Reinforcement Learning(MARL) problem
  - MARL is when we are considering various intelligent agents interacting with an environment
- Multi Agent scenario
  - Cooperative
  - Competitive
  - Mixture of both

# Reinforcement Learning at a glance



RL feedback loop[1]

# Problem Formulation

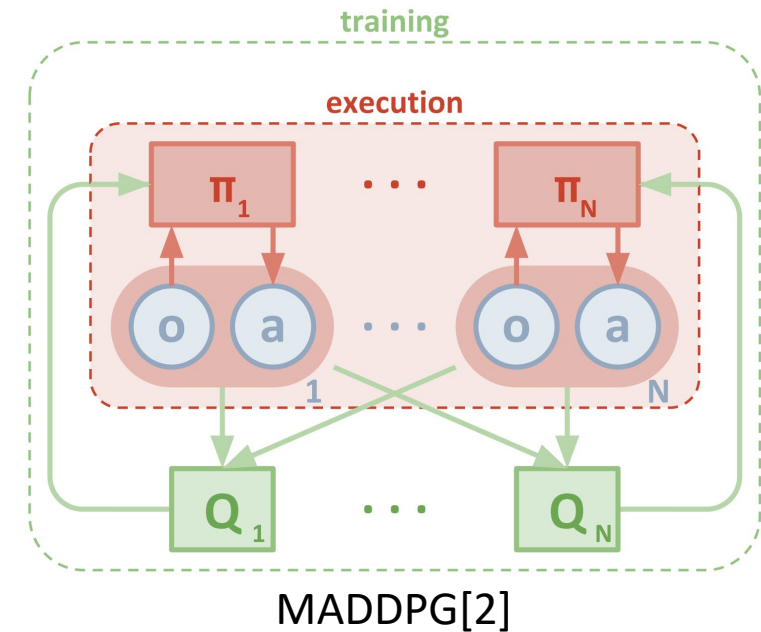
- $N$  number of agents
- Observation space :  $O_t = \{o_1, o_2, o_3, \dots o_N\}$ 
  - $o_i = \{s_1, s_2, \dots s_i, \dots, s_N\}$
- Action space :  $A_t = \{a_1, a_2, a_3, \dots a_N\}$
- Parameterized policy :  $\pi_{\theta_i} : o_i \rightarrow a_i$
- Transition function :  $T : O_t \times A_t \rightarrow O'_{t+1}$
- Reward :  $r_i^t : o_i^t \times a_i^t \rightarrow \mathbb{R}$
- Goal :  $G_i = \sum_t \gamma^t r_i^t$

# Continued

- Actor Critic algorithms
  - Optimizes a parameterized *policy* through policy gradients
- Each agent is modelled as an Actor
  - Updates its policy as suggested by the Critic
- The Critic estimates the values of the state
  - It guides the actor to prefer good actions over bad ones

# MADDPG (Lowe et al)[2]

- Multi Agent Deep Deterministic Policy Gradient
  - Centralized Critic network, Decentralized execution (multiple independent Actors), Deterministic Policy function
- Architecture
  - An actor network for using observations for deterministic actions
  - An identical target actor network for training stability
  - A critic network that uses joint states action pairs to estimate Q-values



# Continued

- Update of the Critic network through loss
  - $L(\theta_c) = \frac{1}{M} \sum_j \left( y_i^j - Q_i^\mu \left( o_i^j, a_1^j, \dots, a_N^j \right) \right)^2$ 
    - where  $y_i^j = r_i^j + \gamma Q_i^{\mu'} \left( o_i'^j, a_1', \dots, a_N' \right) |_{a_i' = \mu_i'(o_i'^j)}$
- Update of the Actor network through the sampled policy gradient
  - $\nabla_{\theta_i} J \approx \frac{1}{M} \sum_j \nabla_{\theta_i} \mu_i \left( o_i^j \right) \nabla_{a_i} Q_i^\mu \left( o_i^j, a_1^j, \dots, a_N^j \right) |_{a_i = \mu_i(o_i^j)}$



# Issues

- No obstacle avoidance : static, dynamic
- No safety guarantees
- Large state space configuration
- Techniques
  - Safety Controllers[3]
  - Control Barrier functions[4]

# Safe Reinforcement Learning<sup>[4]</sup>

- Learning optimal policies while satisfying Safety Constraints
- Safe RL will help in transitioning the control to physical systems
- A major challenge in using reinforcement learning is safety - control policies learned using reinforcement learning typically do not provide any safety guarantees

# Safety Layer<sub>(Dalal et al)[3]</sub>

- Constrained policy optimization
  - $\max_{\theta} E[\sum_t \gamma^t R(o_i, \mu_{i|\theta}(o_i))]$  where  $c_k(o_i) \leq C_k, \forall k \in K$
- Safety Layer
  - $g(s, w)$
- Learning Safety Layer
  - $a_i^* = \arg \min_{a|i} \frac{1}{2} \| a_i - \mu_{i|\theta}(o) \|^2$  where  $c_k(o_i) + g(o, w)^T a_i \leq C_k, \forall k \in K$

# Safe-MADDPG

- Step 1

---

**Algorithm 1:** Safe-MADDPG for N agents

---

```
1: while epochs  $\leq$  TOTAL EPOCHS do
2:   while step  $\leq$  TOTAL SAMPLE STEPS do
3:     Reset Environment(state =  $o_i, \forall i \in N$ );
4:     Calculate constraint values  $c_i$ ;
5:     Execute actions  $(a_1, \dots, a_N)$  and
       observe reward  $r_i$  and new obs  $o'_i$ 
       for each agent  $\forall i \in N$ ;
6:     Calculate constraint values  $c_i^{next}$ ;
7:     Store  $(a_i, o_i, c_i, c_i^{next})$  in replay buffer  $D$ ;
8:   end while
```

- Safety Network update

```
9:   while update via mini-batch do
10:     Sample a mini-batch for each agent  $i \in N$  :
         $(a_i, o_i, c_i, c_i^{next})$ ;
11:     Observe safety layer :  $g_i(o_i, w_i)$ ;
12:     Predict next constraint values;
13:      $c_i^{next.predicted} = c_i + g_i(o_i, w_i)^T a_i$ ;
14:     Calculate loss and update safety layer weights;
15:      $L^k = \| c_i^{next} - c_i^{next.predicted} \|^2, \forall k \in K$ 
16:   end while
17: end while
```

# Safe-MADDPG

## • Step 2

```
18: while episode number  $\leq$  TOTAL EPISODES do
19:   Initialize exploration noise;
20:   Reset environment (state =  $o_i$ ,  $\forall i \in N$ );
21:   while timestep(t)  $\leq$  EPISODE LENGTH do
22:     For an agent i, select action  $a_i = \mu_{i|\theta}(o_i) + \mathcal{N}_t$ 
        w.r.t. the current policy and exploration;
23:     calculate lagrange multipliers :
         $\lambda_i^* = \frac{g(o_i, w_i)^T \mu_{i|\theta}(o_i) + c_i(o_i) - C_i}{g(o_i, w_i)^T g(o_i, w_i)}$ ;
24:     correct each agent's action as :
         $a_i^* = \mu_{i|\theta}(o_i) - \lambda_i^* g(o_i, w_i)$ ;
25:     Execute corrected actions ( $a_1^*, \dots, a_N^*$ ) and
        observe reward  $r_i$  and new obs  $o'_i$  for agent  $i$ ;
26:     Store ( $o_i, a_i^*, r_i, o'_i$ ) in replay buffer  $D$ ;
27:      $o_i \leftarrow o'_i$ ;
```

## • Next

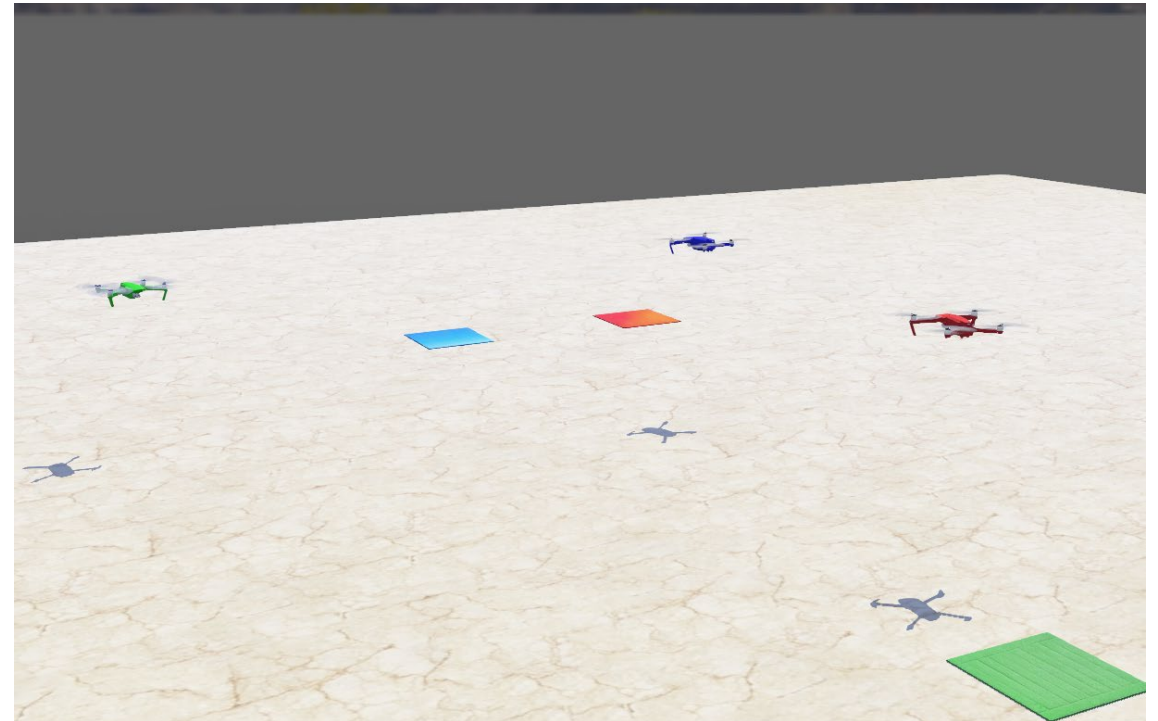
```
28:   while agent i = 1 to N do
29:     Sampling of a randomized minibatch of  $M$ 
        samples ( $o_i^j, a_i^j, r_i^j, o_i'^j$ ) from  $D$ ;
30:     set  $y_i^j =$ 
         $r_i^j + \gamma Q_i^{\mu'}(o_i'^j, a_1', \dots, a_N')|_{a_i' = \mu_{i|\theta'}(o_i'^j)}$ ;
31:     Update of the critic through loss  $L(\theta_c) =$ 
         $\frac{1}{M} \sum_j (y_i^j - Q_i^\mu(o_i^j, a_1^j, \dots, a_N^j))^2$ ;
32:     Update of the actor by the use of sampled
        policy gradient;
         $\nabla_{\theta_i} J \approx \frac{1}{M} \sum_j \nabla_{\theta_i} \mu_i(o_i^j) \nabla_{a_i} Q_i^\mu(o_i^j, a_1^j, \dots, a_N^j)$ 
33:   end while
34:   Update of the target network parameters
        for each agent i;
35:      $\theta'_i \leftarrow \tau \theta_i + (1 - \tau) \theta'_i$ ;
36:   end while
37: end while
```

---

# Experimentation

- Multiple UAVs and their Target locations
- UAV : DJI Mavic
- Webots Simulation by Cyberbotics[5]
  - Real world physics engine
  - Deepbots[6]: Deepbots is a simple framework which is used as middleware between the Webots Simulator and Reinforcement Learning algorithm

- Emitter-Receiver scheme
- Combined robot-supervisor controller scheme



# System Settings

- UAV Observation Space

- Position of each entity (UAV, Target location) :  $P_i : [p_i^x, p_i^y, p_i^z]$
- Velocity of each UAV :  $V_i : [l_i^x, l_i^y, l_i^z, v_i^x, v_i^y, v_i^z]$
- Concatenated tuple :  $\{P^i, V^i, P_{Target}^i, \bigcup_{\forall j \in N; j \neq i} P^j\}$

- UAV Action Space

- For each UAV :  $a_i : \{pitch_i, yaw_i\}$
- Constraints (K=4) :  $\{pitch_{max}, pitch_{min}, yaw_{max}, yaw_{min}\}$

# Continued

- Reward Structure

Reward Value	Conditions
-4	if the angle between agent and its target is $> 1.5$ Radians
-2	if angle $< 1.5$ Radians and $curr\_dist > prev\_dist$
-2	if angle $< 0.1$ Radians and $curr\_dist > prev\_dist$
+5	if angle $< 1.5$ Radians and $curr\_dist < prev\_dist$
+10	if angle $< 0.1$ Radians and $curr\_dist < prev\_dist$
-10	if there is a collision

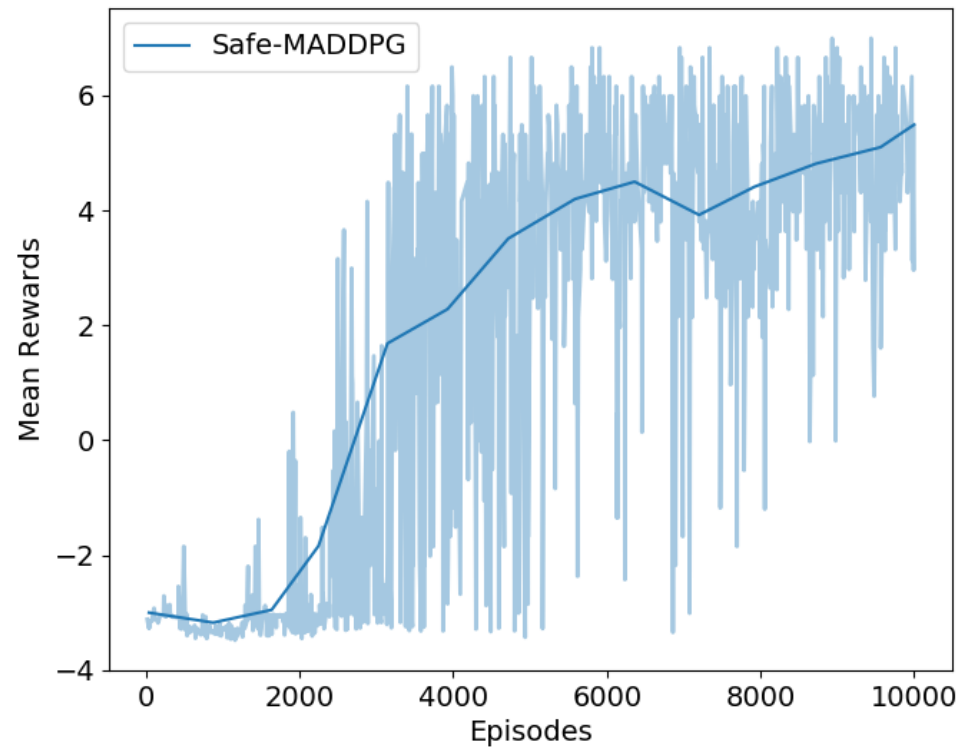


# Results (3 UAV scenario)

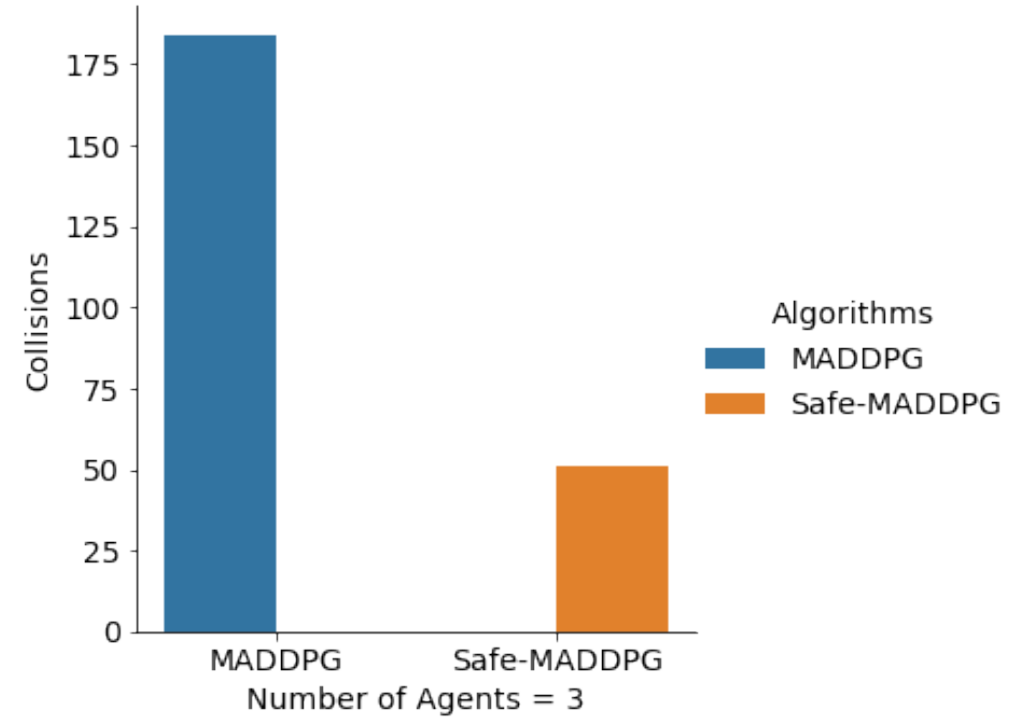
- <video>

# Evaluation

- Mean Rewards



- Testing final policies



# Behavior MADDPG

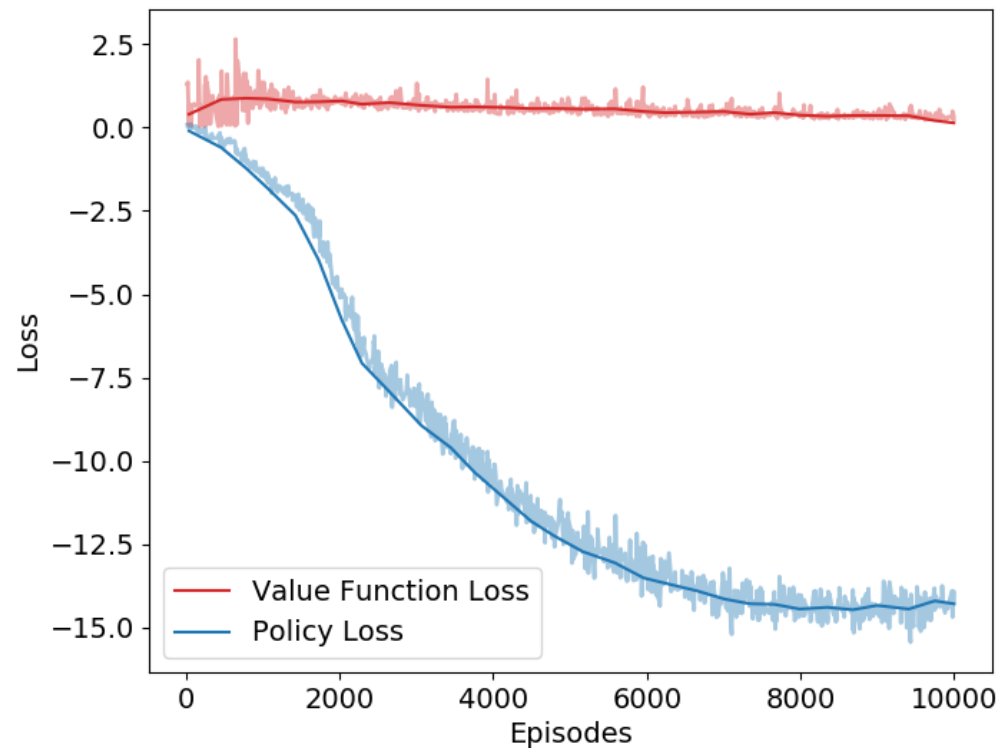
- <video>

# Behavior Safe-MADDPG

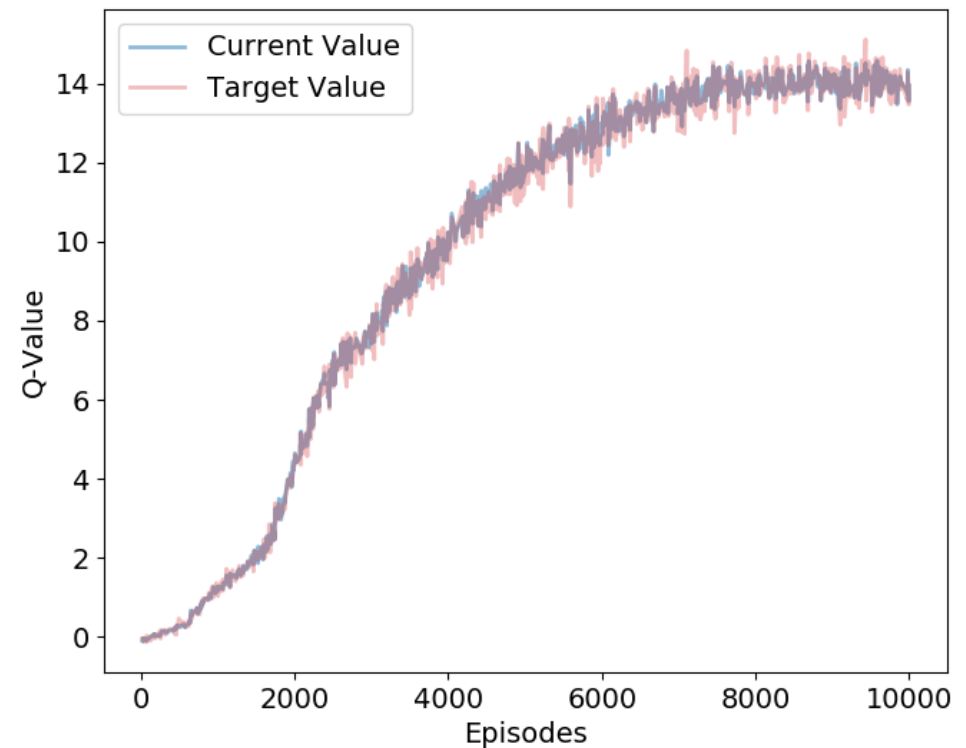
- <video>

# Evaluation

- Value & Policy Loss

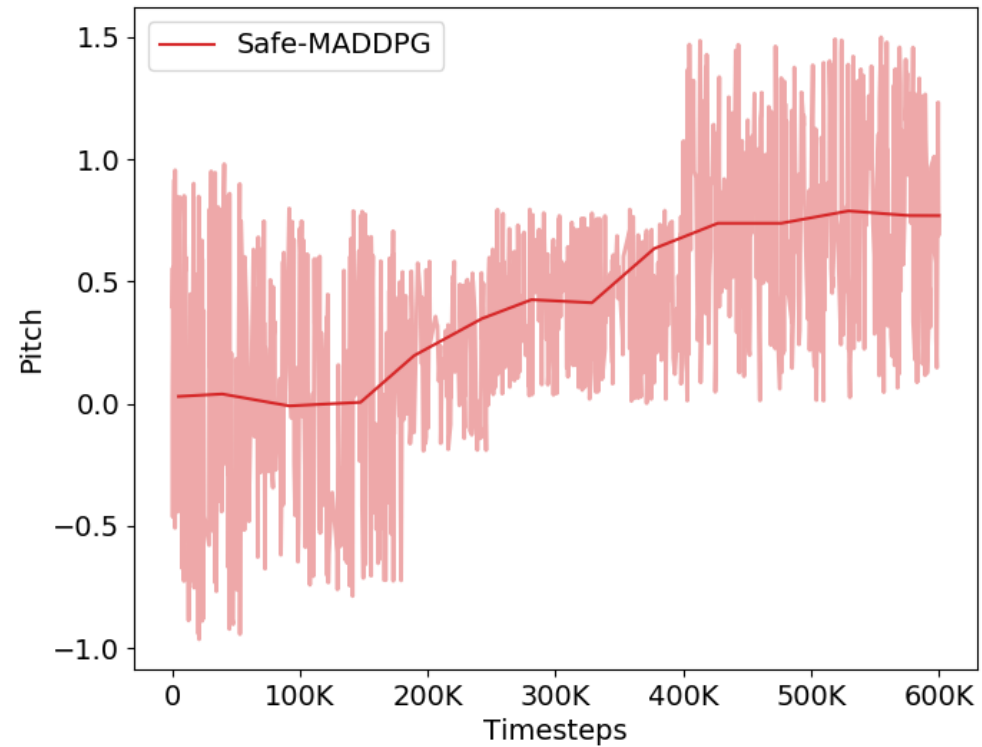


- Current – Target Q-Value

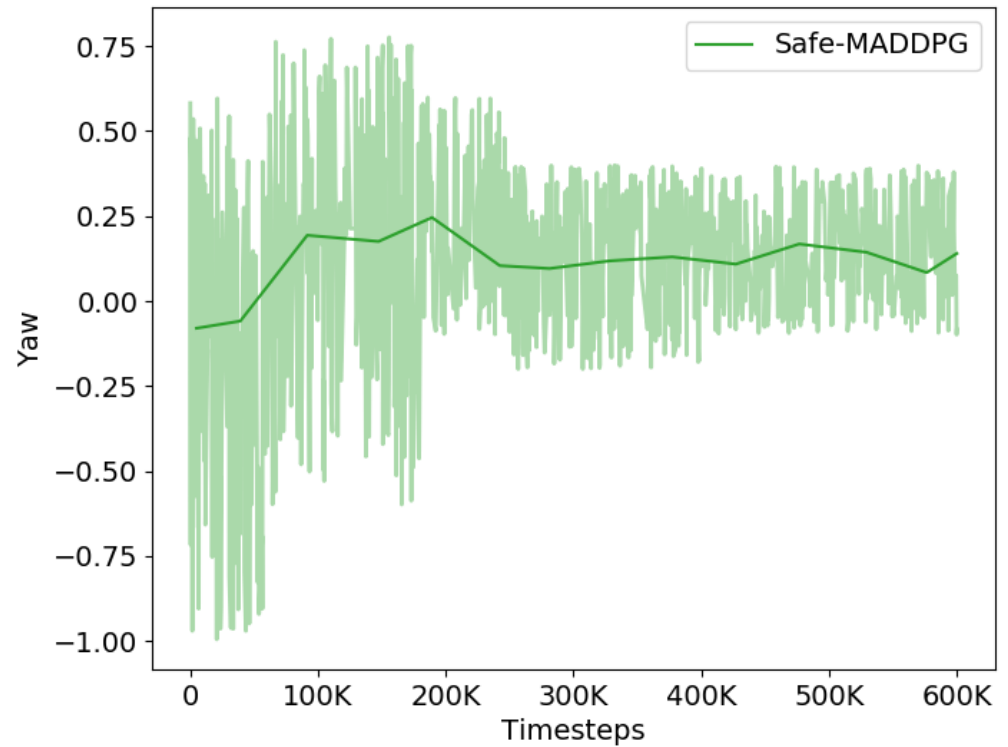


# Evaluation

- Pitch variation



- Yaw variation



# References

1. [Deep Reinforcement Learning \(researchgate.net\)](https://www.researchgate.net)
2. R.Lowe, Y.Wu, A.Tamar, J.Harb, P.Abbeel, and I. Mordatch, “Multi-agent actor-critic for mixed cooperative-competitive environments,” *CoRR*, vol. abs/1706.02275, 2017.
3. G.Dalal, K.Dvijotham, M. Vecerik, T. Hester, C.Paduraru, and Y.Tassa, “Safe exploration in continuous action spaces,” *CoRR*, vol. abs/1801.08757, 2018.
4. R.Cheng, G.Orosz, R.M.Murray, and J.W.Burdick, “End-to-end safe reinforcement learning through barrier functions for safety-critical continuous control tasks,” in *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 33, 2019.
5. Webots, “<http://www.cyberbotics.com>,” open-source Mobile Robot Simulation Software. [Online]. Available: <http://www.cyberbotics.com>.
6. M.Kirtas, K.Tsampazis, N.Passalis, and A.Tefas, “Deepbots: A webots-based deep reinforcement learning framework for robotics,” in *Artificial Intelligence Applications and Innovations*. Cham: Springer International Publishing, 2020.