

Communicating Multi-agent Collision Avoidance with Deep Reinforcement Learning

By

Yu Fan Chen, Miao Liu, Michael Everett, and Jonathan P.

How

Presenter: Jared Choi

Motivation

- Finding a path
- Computationally expensive due to
 - Collision checking
 - Feasibility checking
 - Efficiency checking



Motivation

- Finding a path
 - Computationally expensive due to
 - Collision checking
 - Feasibility checking
 - Efficiency checking
- Offline Learning



Background

- A sequential decision making problem can be formulated as a Markov Decision Process (MDP)
- $M = \langle S, A, P, \gamma, R, \quad \rangle$

Background

- A sequential decision making problem can be formulated as a Markov Decision Process (MDP)
- $M = \langle S, A, P, \gamma, R \rangle$
 - S (state space)
 - A (action space)
 - P (state transition model)
 - R : reward function
 - γ : discount factor

State Space ($M = \langle \mathbf{S}, \mathcal{A}, P, R, \gamma \rangle$)

- S (state space)
 - System's state is constructed by concatenating the two agents' individual states

$$s^o = [p_x, p_y, v_x, v_y, r] \in \mathbb{R}^5$$

Observable State Vector
(position (x,y), velocity(x,y), radius)

$$s^h = [p_{gx}, p_{gy}, v_{pref}, \psi] \in \mathbb{R}^4$$

Unobservable State vector
(goal position (x,y), preferred speed, heading)

State Space ($M = \langle \mathbf{S}, \mathcal{A}, P, R, \gamma \rangle$)

- \mathbf{S} (state space)

- System's state is constructed by concatenating the two agents' individual states

$$s^o = [p_x, p_y, v_x, v_y, r] \in \mathbb{R}^5 \quad \text{Observable State Vector (position (x,y), velocity(x,y), radius)}$$

$$s^h = [p_{gx}, p_{gy}, v_{pref}, \psi] \in \mathbb{R}^4 \quad \text{Unobservable State vector (goal position (x,y), preferred speed, heading)}$$

$$\mathbf{s}^{jn} = [\mathbf{s}, \tilde{\mathbf{s}}^o] \in \mathbb{R}^{14}$$

Action Space ($M = \langle S, \mathbf{A}, P, R, \gamma \rangle$)

- A(action space):
 - Set of permissible velocity vectors, $a(s) = \mathbf{v}$

$$\mathbf{a}(s) = \mathbf{v} \text{ for } \|\mathbf{v}\|_2 < v_{pref}$$

State Transition Model ($M = \langle S, \gamma, A, \mathbf{P}, R, \gamma \rangle$)

- P(state transition model)
 - A probabilistic state transition model
- Determined by the agents' kinematics
- Unknown to us $P(\mathbf{s}_{t+1}^{jn}, \mathbf{s}_t^{j'n} | \mathbf{a}_t)$

Reward Function ($M = \langle S, A, P, \mathbf{R}, \gamma \rangle$)

- R : reward function
 - Award the agent for reaching its goal
 - Penalize the agent for getting too close or colliding with other agent

$$R(\mathbf{s}^{jn}, \mathbf{a}) = \begin{cases} -0.25 & \text{if } d_{min} < 0 \\ -0.1 - d_{min}/2 & \text{else if } d_{min} < 0.2 \\ 1 & \text{else if } \mathbf{p} = \mathbf{p}_g \\ 0 & \text{o.w.} \end{cases}$$

Discount Factor($M = \langle S, A, P, R, \gamma \rangle$)

- Discount factor

$$\gamma \in [0, 1)$$

Value Function

- The value of a state
- Value depends on
 - γ close to 1
 - We care about our long term reward
 - γ close to 0
 - We care only about our immediate reward

$$V^*(\mathbf{s}_0^{jn}) = \sum_{t=0}^T \gamma^{t \cdot v_{pref}} R(\mathbf{s}_t^{jn}, \pi^*(\mathbf{s}_t^{jn}))$$

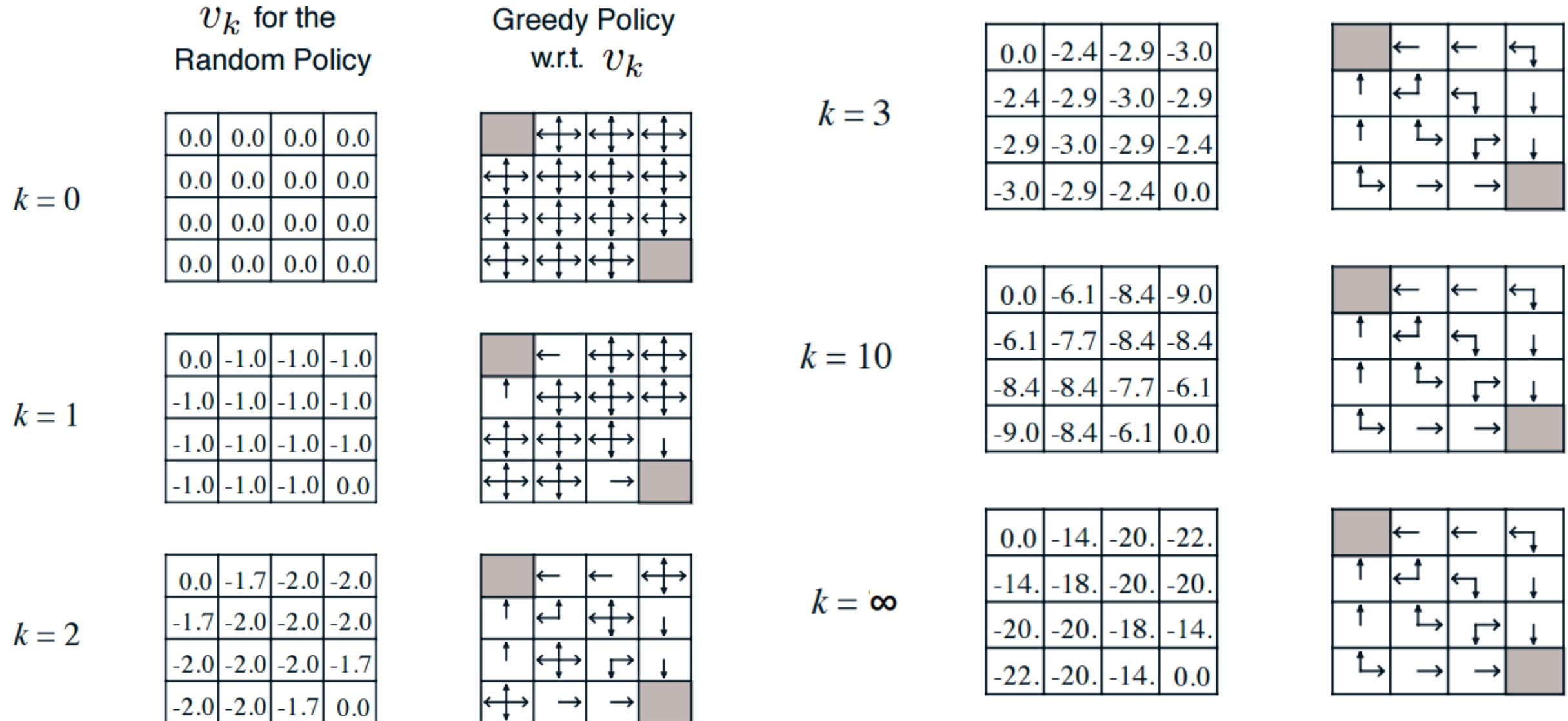
Optimal Policy

- The best trajectory at given state

$$\pi^*(\mathbf{s}_0^{jn}) = \operatorname{argmax}_{\mathbf{a}} R(\mathbf{s}_0, \mathbf{a}) + \gamma^{\Delta t \cdot v_{pref}} \int_{\mathbf{s}_1^{jn}} P(\mathbf{s}_0^{jn}, \mathbf{s}_1^{jn} | \mathbf{a}) V^*(\mathbf{s}_1^{jn}) d\mathbf{s}_1^{jn}$$

Value Function and Optimal Policy

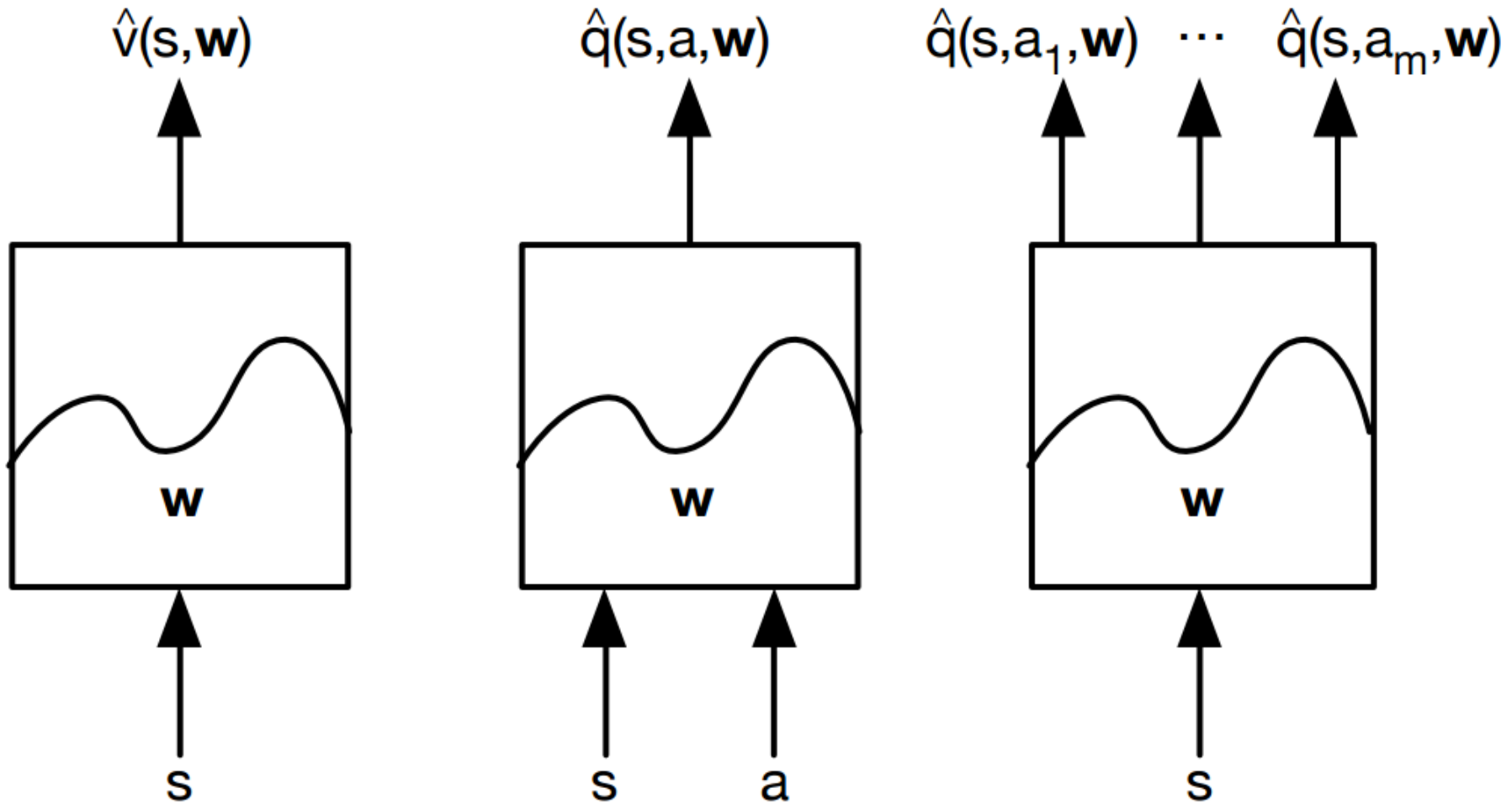
From David Silver's slides



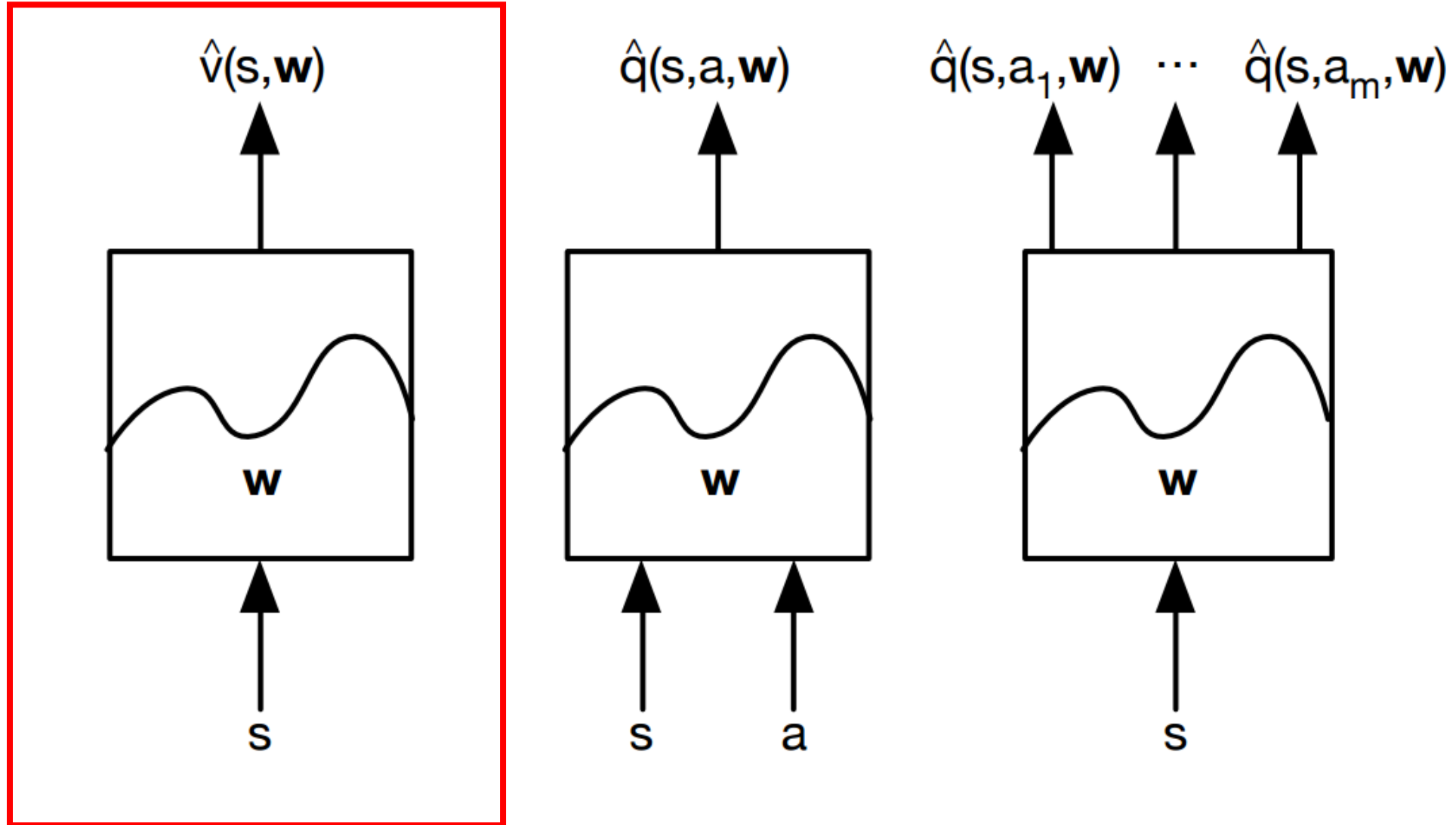
Value Function and Optimal Policy

- Every state s has value $V(s)$
 - Store it in a lookup table
 - In a grid world : 16 values
 - In motion planning : Infinite values (b/c it's continuous state space)
- Solution:
 - Approximate value via neural network

Value Function and Optimal Policy



Value Function and Optimal Policy



Collision Avoidance Deep Reinforcement Learning

1. Train Value network using ORCA

2. Train again with Deep reinforcement Learning

Collision Avoidance Deep Reinforcement Learning

1. Train Value network using ORCA

- Why pre-train?

Collision Avoidance Deep Reinforcement Learning

1. Train Value network using ORCA

- Why pre-train?
 - Initializing the neural network is crucial to convergence
 - We want the network to output something reasonable

Collision Avoidance Deep Reinforcement Learning

1. Train Value network using ORCA

- Why pre-train?
 - Initializing the neural network is crucial to convergence
 - We want the network to output something reasonable
- Generate 500 trajectories as a training set
 - Each trajectory contains 40 state-value pairs (total of 20,000 pairs)
 - Back-propagate to minimize our loss function.

$$\operatorname{argmin}_{\mathbf{w}} \sum_{k=1}^N \left(y_k - V(\mathbf{s}_k^{jn}; \mathbf{w}) \right)^2$$

Collision Avoidance Deep Reinforcement Learning

1. Train Value network using ORCA

2. Train again with Deep reinforcement Learning

Collision Avoidance Deep Reinforcement Learning

1. Train again with Deep reinforcement Learning

Algorithm 2: Deep V-learning

```

1 Input: trajectory training set  $D$ 
2 Output: value network  $V(\cdot; \mathbf{w})$ 
3  $V(\cdot; \mathbf{w}) \leftarrow \text{train\_nn}(D)$  //step 1: initialization
4 duplicate value net  $V' \leftarrow V$  //step 2: RL
5 initialize experience set  $E \leftarrow D$ 
6 for  $episode=1, \dots, N_{eps}$  do
7   for  $m$  times do
8      $\mathbf{s}_0, \tilde{\mathbf{s}}_0 \leftarrow \text{randomTestcase}()$ 
9      $\mathbf{s}_{0:t_f} \leftarrow \text{CADRL}(V), \tilde{\mathbf{s}}_{0:\tilde{t}_f} \leftarrow \text{CADRL}(V)$ 
10     $\mathbf{y}_{0:T}, \tilde{\mathbf{y}}_{0:\tilde{t}_f} \leftarrow \text{findValues}(V', \mathbf{s}_{0:t_f}, \tilde{\mathbf{s}}_{0:\tilde{t}_f})$ 
11     $E \leftarrow \text{assimilate}(E, (\mathbf{y}, \mathbf{s}^{jn})_{0:t_f}, (\tilde{\mathbf{y}}, \tilde{\mathbf{s}}^{jn})_{0:\tilde{t}_f})$ 
12     $e \leftarrow \text{randSubset}(E)$ 
13     $\mathbf{w} \leftarrow \text{backprop}(e)$ 
14    for every  $C$  episodes do
15      Evaluate( $V$ ),  $V' \leftarrow V$ 
16 return  $V$ 

```

Collision Avoidance Deep Reinforcement Learning

1. Train again with Deep reinforcement Learning

Algorithm 2: Deep V-learning

```

1 Input: trajectory training set  $D$ 
2 Output: value network  $V(\cdot; \mathbf{w})$ 
3  $V(\cdot; \mathbf{w}) \leftarrow \text{train\_nn}(D)$  //step 1: initialization
4 duplicate value net  $V' \leftarrow V$  //step 2: RL
5 initialize experience set  $E \leftarrow D$ 
6 for  $episode=1, \dots, N_{eps}$  do
7   for  $m$  times do
8      $\mathbf{s}_0, \tilde{\mathbf{s}}_0 \leftarrow \text{randomTestcase}()$ 
9      $\mathbf{s}_{0:t_f} \leftarrow \text{CADRL}(V), \tilde{\mathbf{s}}_{0:\tilde{t}_f} \leftarrow \text{CADRL}(V)$ 
10     $\mathbf{y}_{0:T}, \tilde{\mathbf{y}}_{0:\tilde{t}_f} \leftarrow \text{findValues}(V', \mathbf{s}_{0:t_f}, \tilde{\mathbf{s}}_{0:\tilde{t}_f})$ 
11     $E \leftarrow \text{assimilate}(E, (\mathbf{y}, \mathbf{s}^{jn})_{0:t_f}, (\tilde{\mathbf{y}}, \tilde{\mathbf{s}}^{jn})_{0:\tilde{t}_f})$ 
12     $e \leftarrow \text{randSubset}(E)$ 
13     $\mathbf{w} \leftarrow \text{backprop}(e)$ 
14    for every  $C$  episodes do
15      Evaluate( $V$ ),  $V' \leftarrow V$ 
16 return  $V$ 

```


Collision Avoidance Deep Reinforcement Learning

1. Train again with Deep reinforcement Learning

Algorithm 2: Deep V-learning

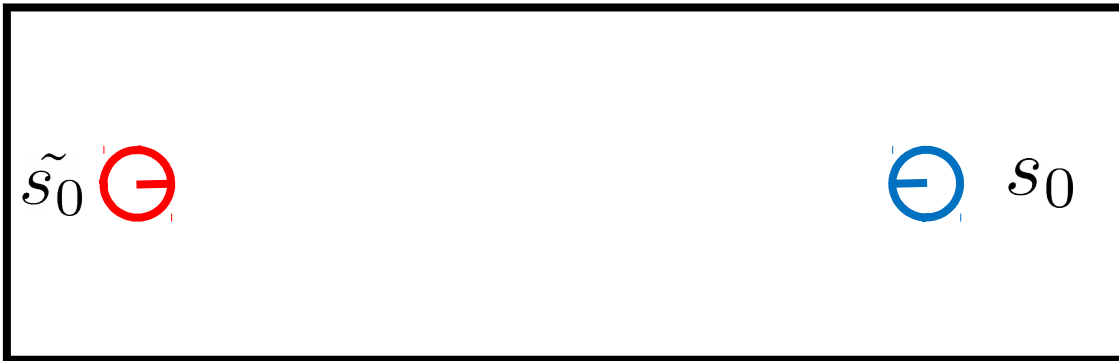
```

1 Input: trajectory training set  $D$ 
2 Output: value network  $V(\cdot; \mathbf{w})$ 
3  $V(\cdot; \mathbf{w}) \leftarrow \text{train\_nn}(D)$  //step 1: initialization
4 duplicate value net  $V' \leftarrow V$  //step 2: RL
5 initialize experience set  $E \leftarrow D$ 
6 for  $episode=1, \dots, N_{eps}$  do
7   for  $m$  times do
8      $\mathbf{s}_0, \tilde{\mathbf{s}}_0 \leftarrow \text{randomTestcase}()$ 
9      $\mathbf{s}_{0:t_f} \leftarrow \text{CADRL}(V), \tilde{\mathbf{s}}_{0:\tilde{t}_f} \leftarrow \text{CADRL}(V)$ 
10     $\mathbf{y}_{0:T}, \tilde{\mathbf{y}}_{0:\tilde{t}_f} \leftarrow \text{findValues}(V', \mathbf{s}_{0:t_f}, \tilde{\mathbf{s}}_{0:\tilde{t}_f})$ 
11     $E \leftarrow \text{assimilate}(E, (\mathbf{y}, \mathbf{s}^{jn})_{0:t_f}, (\tilde{\mathbf{y}}, \tilde{\mathbf{s}}^{jn})_{0:\tilde{t}_f})$ 
12     $e \leftarrow \text{randSubset}(E)$ 
13     $\mathbf{w} \leftarrow \text{backprop}(e)$ 
14    for every  $C$  episodes do
15      Evaluate( $V$ ),  $V' \leftarrow V$ 
16 return  $V$ 

```

Collision Avoidance Deep Reinforcement Learning

1. Train again with Deep reinforcement Learning



Algorithm 2: Deep V-learning

```

1 Input: trajectory training set  $D$ 
2 Output: value network  $V(\cdot; \mathbf{w})$ 
3  $V(\cdot; \mathbf{w}) \leftarrow \text{train\_nn}(D)$  //step 1: initialization
4 duplicate value net  $V' \leftarrow V$  //step 2: RL
5 initialize experience set  $E \leftarrow D$ 
6 for  $episode=1, \dots, N_{eps}$  do
7   for  $m$  times do
8      $s_0, \tilde{s}_0 \leftarrow \text{randomTestcase}()$ 
9      $s_{0:t_f} \leftarrow \text{CADRL}(V), \tilde{s}_{0:\tilde{t}_f} \leftarrow \text{CADRL}(V)$ 
10     $\mathbf{y}_{0:T}, \tilde{\mathbf{y}}_{0:\tilde{t}_f} \leftarrow \text{findValues}(V', s_{0:t_f}, \tilde{s}_{0:\tilde{t}_f})$ 
11     $E \leftarrow \text{assimilate}(E, (\mathbf{y}, \mathbf{s}^{jn})_{0:t_f}, (\tilde{\mathbf{y}}, \tilde{\mathbf{s}}^{jn})_{0:\tilde{t}_f})$ 
12     $e \leftarrow \text{randSubset}(E)$ 
13     $\mathbf{w} \leftarrow \text{backprop}(e)$ 
14    for every  $C$  episodes do
15      Evaluate( $V$ ),  $V' \leftarrow V$ 
16 return  $V$ 

```

Collision Avoidance Deep Reinforcement Learning

1. Train again with Deep reinforcement Learning



Algorithm 2: Deep V-learning

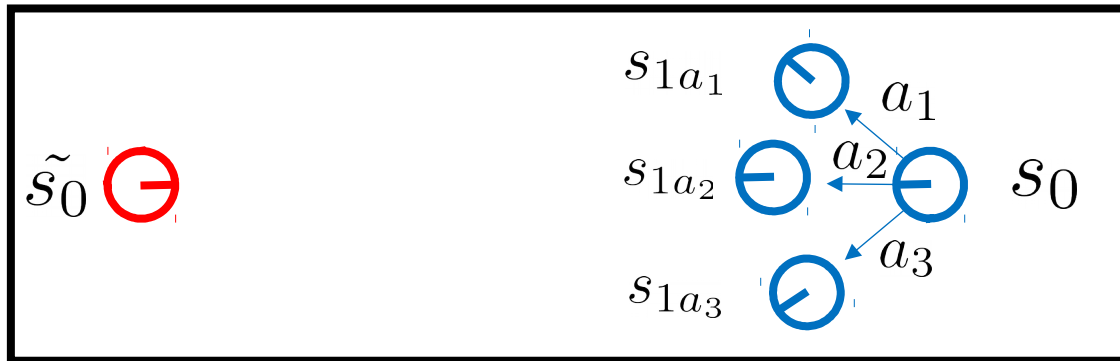
```

1 Input: trajectory training set  $D$ 
2 Output: value network  $V(\cdot; \mathbf{w})$ 
3  $V(\cdot; \mathbf{w}) \leftarrow \text{train\_nn}(D)$  //step 1: initialization
4 duplicate value net  $V' \leftarrow V$  //step 2: RL
5 initialize experience set  $E \leftarrow D$ 
6 for  $episode=1, \dots, N_{eps}$  do
7   for  $m$  times do
8      $s_0, \tilde{s}_0 \leftarrow \text{randomTestcase}()$ 
9      $s_{0:t_f} \leftarrow \text{CADRL}(V), \tilde{s}_{0:\tilde{t}_f} \leftarrow \text{CADRL}(V)$ 
10     $\mathbf{y}_{0:T}, \tilde{\mathbf{y}}_{0:\tilde{t}_f} \leftarrow \text{findValues}(V', s_{0:t_f}, \tilde{s}_{0:\tilde{t}_f})$ 
11     $E \leftarrow \text{assimilate}(E, (\mathbf{y}, \mathbf{s}^{jn})_{0:t_f}, (\tilde{\mathbf{y}}, \tilde{\mathbf{s}}^{jn})_{0:\tilde{t}_f})$ 
12     $e \leftarrow \text{randSubset}(E)$ 
13     $\mathbf{w} \leftarrow \text{backprop}(e)$ 
14    for every  $C$  episodes do
15      Evaluate( $V$ ),  $V' \leftarrow V$ 
16 return  $V$ 

```

Collision Avoidance Deep Reinforcement Learning

1. Train again with Deep reinforcement Learning



Algorithm 2: Deep V-learning

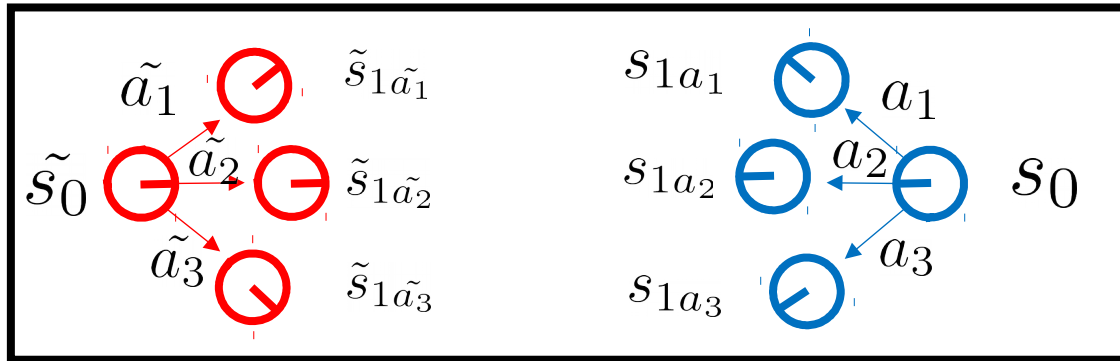
```

1 Input: trajectory training set  $D$ 
2 Output: value network  $V(\cdot; \mathbf{w})$ 
3  $V(\cdot; \mathbf{w}) \leftarrow \text{train\_nn}(D)$  //step 1: initialization
4 duplicate value net  $V' \leftarrow V$  //step 2: RL
5 initialize experience set  $E \leftarrow D$ 
6 for  $episode=1, \dots, N_{eps}$  do
7   for  $m$  times do
8      $s_0, \tilde{s}_0 \leftarrow \text{randomTestcase}()$ 
9      $s_{0:t_f} \leftarrow \text{CADRL}(V), \tilde{s}_{0:\tilde{t}_f} \leftarrow \text{CADRL}(V)$ 
10     $\mathbf{y}_{0:T}, \tilde{\mathbf{y}}_{0:\tilde{t}_f} \leftarrow \text{findValues}(V', s_{0:t_f}, \tilde{s}_{0:\tilde{t}_f})$ 
11     $E \leftarrow \text{assimilate}(E, (\mathbf{y}, \mathbf{s}^{jn})_{0:t_f}, (\tilde{\mathbf{y}}, \tilde{\mathbf{s}}^{jn})_{0:\tilde{t}_f})$ 
12     $e \leftarrow \text{randSubset}(E)$ 
13     $\mathbf{w} \leftarrow \text{backprop}(e)$ 
14    for every  $C$  episodes do
15      Evaluate( $V$ ),  $V' \leftarrow V$ 
16 return  $V$ 

```

Collision Avoidance Deep Reinforcement Learning

1. Train again with Deep reinforcement Learning



Algorithm 2: Deep V-learning

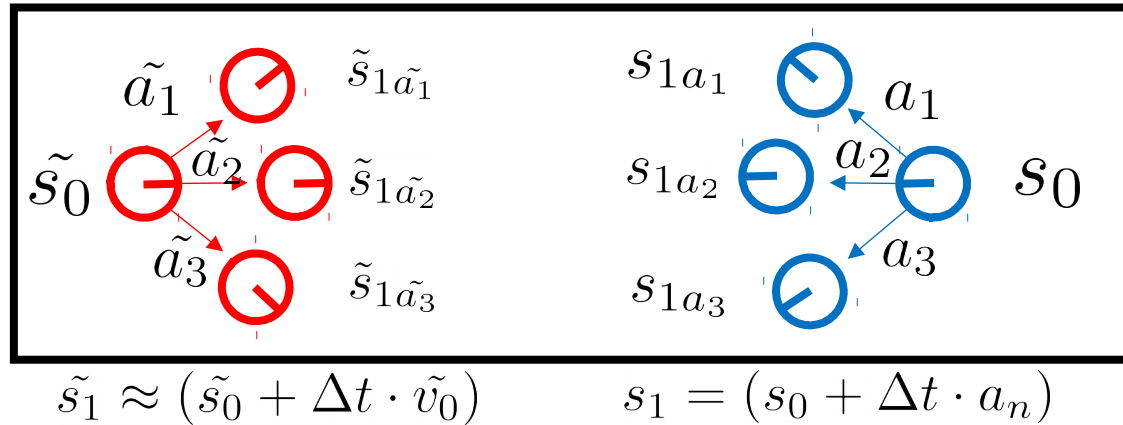
```

1 Input: trajectory training set  $D$ 
2 Output: value network  $V(\cdot; \mathbf{w})$ 
3  $V(\cdot; \mathbf{w}) \leftarrow \text{train\_nn}(D)$  //step 1: initialization
4 duplicate value net  $V' \leftarrow V$  //step 2: RL
5 initialize experience set  $E \leftarrow D$ 
6 for  $episode=1, \dots, N_{eps}$  do
7   for  $m$  times do
8      $s_0, \tilde{s}_0 \leftarrow \text{randomTestcase}()$ 
9      $s_{0:t_f} \leftarrow \text{CADRL}(V), \tilde{s}_{0:\tilde{t}_f} \leftarrow \text{CADRL}(V)$ 
10     $\mathbf{y}_{0:T}, \tilde{\mathbf{y}}_{0:\tilde{t}_f} \leftarrow \text{findValues}(V', s_{0:t_f}, \tilde{s}_{0:\tilde{t}_f})$ 
11     $E \leftarrow \text{assimilate}(E, (\mathbf{y}, \mathbf{s}^{jn})_{0:t_f}, (\tilde{\mathbf{y}}, \tilde{\mathbf{s}}^{jn})_{0:\tilde{t}_f})$ 
12     $e \leftarrow \text{randSubset}(E)$ 
13     $\mathbf{w} \leftarrow \text{backprop}(e)$ 
14    for every  $C$  episodes do
15      Evaluate( $V$ ),  $V' \leftarrow V$ 
16 return  $V$ 

```

Collision Avoidance Deep Reinforcement Learning

1. Train again with Deep reinforcement Learning



Algorithm 2: Deep V-learning

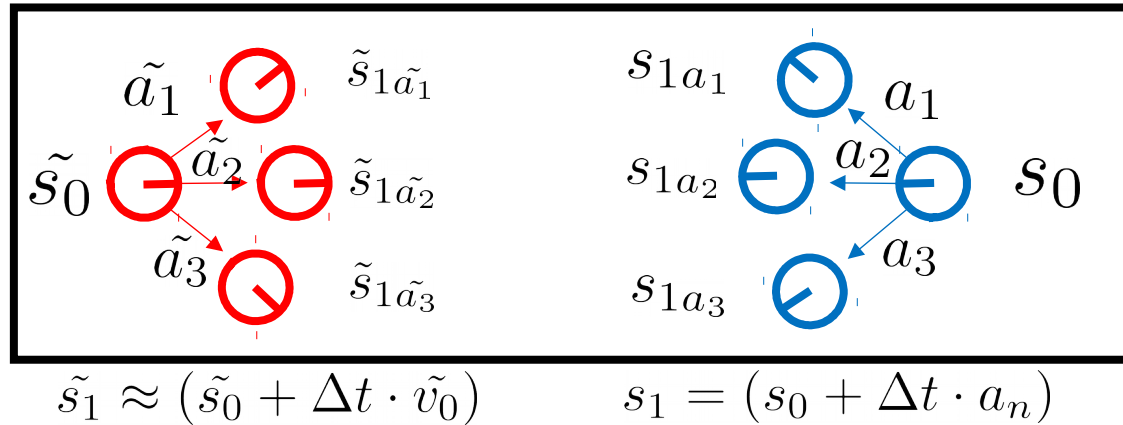
```

1 Input: trajectory training set  $D$ 
2 Output: value network  $V(\cdot; \mathbf{w})$ 
3  $V(\cdot; \mathbf{w}) \leftarrow \text{train\_nn}(D)$  //step 1: initialization
4 duplicate value net  $V' \leftarrow V$  //step 2: RL
5 initialize experience set  $E \leftarrow D$ 
6 for  $episode=1, \dots, N_{eps}$  do
7   for  $m$  times do
8      $s_0, \tilde{s}_0 \leftarrow \text{randomTestcase}()$ 
9      $s_{0:t_f} \leftarrow \text{CADRL}(V), \tilde{s}_{0:\tilde{t}_f} \leftarrow \text{CADRL}(V)$ 
10     $\mathbf{y}_{0:T}, \tilde{\mathbf{y}}_{0:\tilde{t}_f} \leftarrow \text{findValues}(V', s_{0:t_f}, \tilde{s}_{0:\tilde{t}_f})$ 
11     $E \leftarrow \text{assimilate}(E, (\mathbf{y}, \mathbf{s}^{jn})_{0:t_f}, (\tilde{\mathbf{y}}, \tilde{\mathbf{s}}^{jn})_{0:\tilde{t}_f})$ 
12     $e \leftarrow \text{randSubset}(E)$ 
13     $\mathbf{w} \leftarrow \text{backprop}(e)$ 
14    for every  $C$  episodes do
15      Evaluate( $V$ ),  $V' \leftarrow V$ 
16 return  $V$ 

```

Collision Avoidance Deep Reinforcement Learning

1. Train again with Deep reinforcement Learning



$[s_{1a_1}, \tilde{s}_{1\tilde{a}_1}]$

Algorithm 2: Deep V-learning

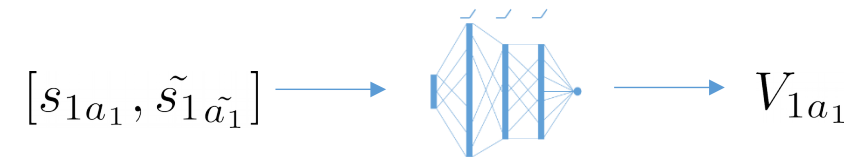
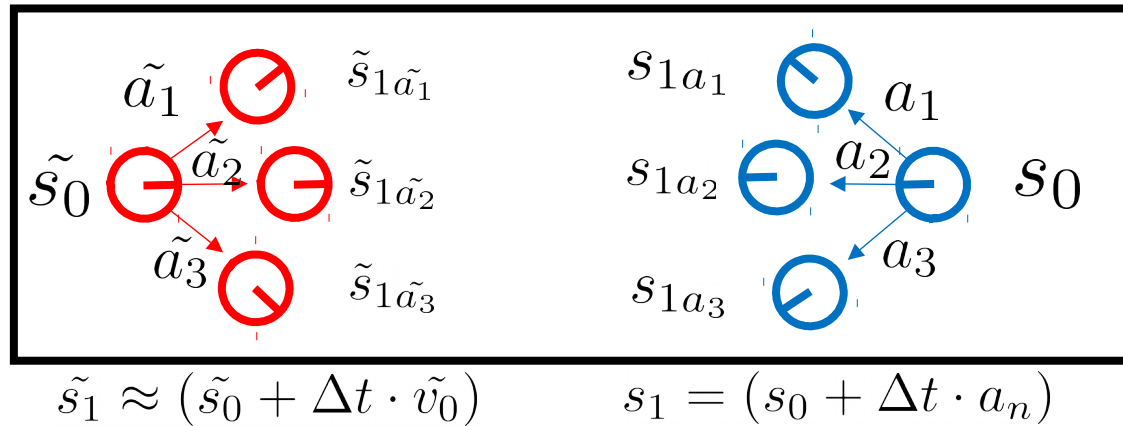
```

1 Input: trajectory training set  $D$ 
2 Output: value network  $V(\cdot; \mathbf{w})$ 
3  $V(\cdot; \mathbf{w}) \leftarrow \text{train\_nn}(D)$  //step 1: initialization
4 duplicate value net  $V' \leftarrow V$  //step 2: RL
5 initialize experience set  $E \leftarrow D$ 
6 for  $episode=1, \dots, N_{eps}$  do
7   for  $m$  times do
8      $s_0, \tilde{s}_0 \leftarrow \text{randomTestcase}()$ 
9      $s_{0:t_f} \leftarrow \text{CADRL}(V), \tilde{s}_{0:\tilde{t}_f} \leftarrow \text{CADRL}(V)$ 
10     $\mathbf{y}_{0:T}, \tilde{\mathbf{y}}_{0:\tilde{t}_f} \leftarrow \text{findValues}(V', s_{0:t_f}, \tilde{s}_{0:\tilde{t}_f})$ 
11     $E \leftarrow \text{assimilate}(E, (\mathbf{y}, \mathbf{s}^{jn})_{0:t_f}, (\tilde{\mathbf{y}}, \tilde{\mathbf{s}}^{jn})_{0:\tilde{t}_f})$ 
12     $e \leftarrow \text{randSubset}(E)$ 
13     $\mathbf{w} \leftarrow \text{backprop}(e)$ 
14    for every  $C$  episodes do
15      Evaluate( $V$ ),  $V' \leftarrow V$ 
16 return  $V$ 

```

Collision Avoidance Deep Reinforcement Learning

1. Train again with Deep reinforcement Learning

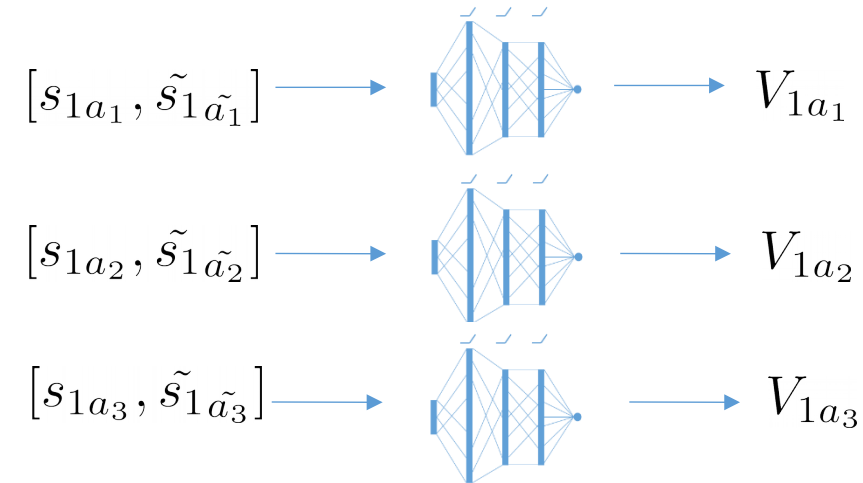
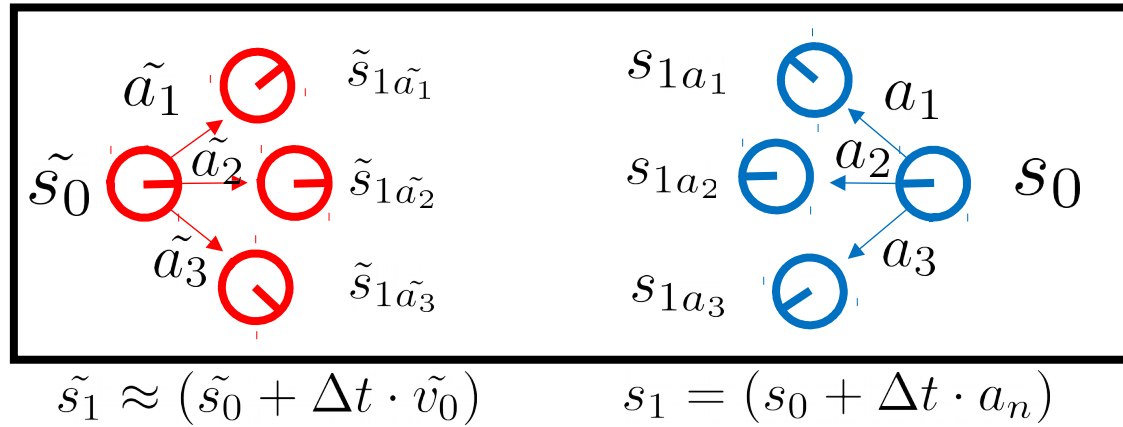


Algorithm 2: Deep V-learning

- 1 **Input:** trajectory training set D
- 2 **Output:** value network $V(\cdot; \mathbf{w})$
- 3 $V(\cdot; \mathbf{w}) \leftarrow \text{train_nn}(D)$ //step 1: initialization
- 4 duplicate value net $V' \leftarrow V$ //step 2: RL
- 5 initialize experience set $E \leftarrow D$
- 6 **for** $episode=1, \dots, N_{eps}$ **do**
- 7 **for** m times **do**
- 8 $s_0, \tilde{s}_0 \leftarrow \text{randomTestcase}()$
- 9 $s_{0:t_f} \leftarrow \text{CADRL}(V), \tilde{s}_{0:\tilde{t}_f} \leftarrow \text{CADRL}(V)$
- 10 $\mathbf{y}_{0:T}, \tilde{\mathbf{y}}_{0:\tilde{t}_f} \leftarrow \text{findValues}(V', s_{0:t_f}, \tilde{s}_{0:\tilde{t}_f})$
- 11 $E \leftarrow \text{assimilate}(E, (\mathbf{y}, \mathbf{s}^{jn})_{0:t_f}, (\tilde{\mathbf{y}}, \tilde{\mathbf{s}}^{jn})_{0:\tilde{t}_f})$
- 12 $e \leftarrow \text{randSubset}(E)$
- 13 $\mathbf{w} \leftarrow \text{backprop}(e)$
- 14 **for every** C episodes **do**
- 15 Evaluate(V), $V' \leftarrow V$
- 16 **return** V

Collision Avoidance Deep Reinforcement Learning

1. Train again with Deep reinforcement Learning

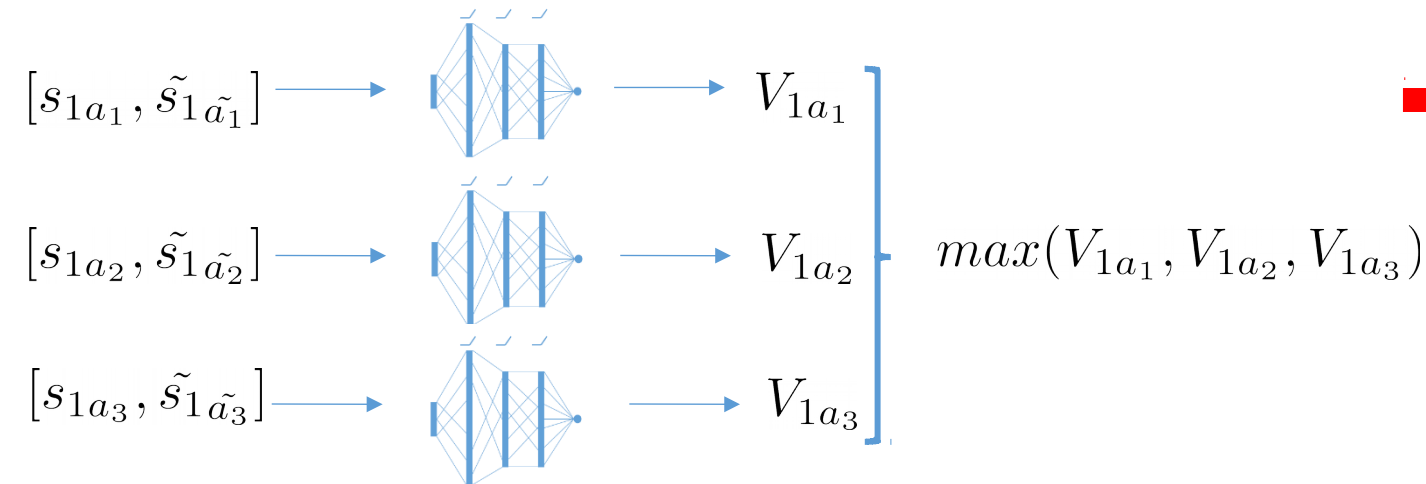
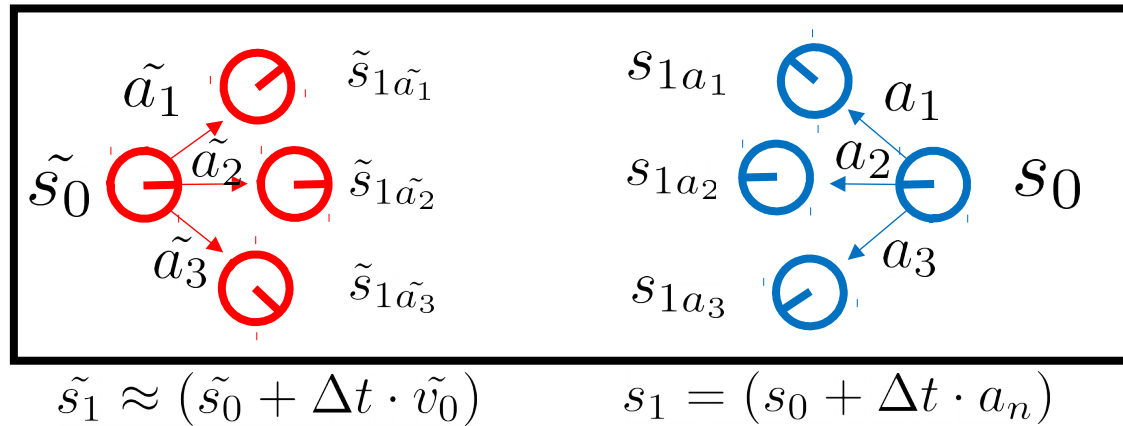


Algorithm 2: Deep V-learning

- 1 **Input:** trajectory training set D
- 2 **Output:** value network $V(\cdot; \mathbf{w})$
- 3 $V(\cdot; \mathbf{w}) \leftarrow \text{train_nn}(D)$ //step 1: initialization
- 4 duplicate value net $V' \leftarrow V$ //step 2: RL
- 5 initialize experience set $E \leftarrow D$
- 6 **for** $episode=1, \dots, N_{eps}$ **do**
- 7 **for** m times **do**
- 8 $s_0, \tilde{s}_0 \leftarrow \text{randomTestcase}()$
- 9 $s_{0:t_f} \leftarrow \text{CADRL}(V), \tilde{s}_{0:\tilde{t}_f} \leftarrow \text{CADRL}(V)$
- 10 $\mathbf{y}_{0:T}, \tilde{\mathbf{y}}_{0:\tilde{t}_f} \leftarrow \text{findValues}(V', s_{0:t_f}, \tilde{s}_{0:\tilde{t}_f})$
- 11 $E \leftarrow \text{assimilate}(E, (\mathbf{y}, \mathbf{s}^{jn})_{0:t_f}, (\tilde{\mathbf{y}}, \tilde{\mathbf{s}}^{jn})_{0:\tilde{t}_f})$
- 12 $e \leftarrow \text{randSubset}(E)$
- 13 $\mathbf{w} \leftarrow \text{backprop}(e)$
- 14 **for** every C episodes **do**
- 15 Evaluate(V), $V' \leftarrow V$
- 16 **return** V

Collision Avoidance Deep Reinforcement Learning

1. Train again with Deep reinforcement Learning

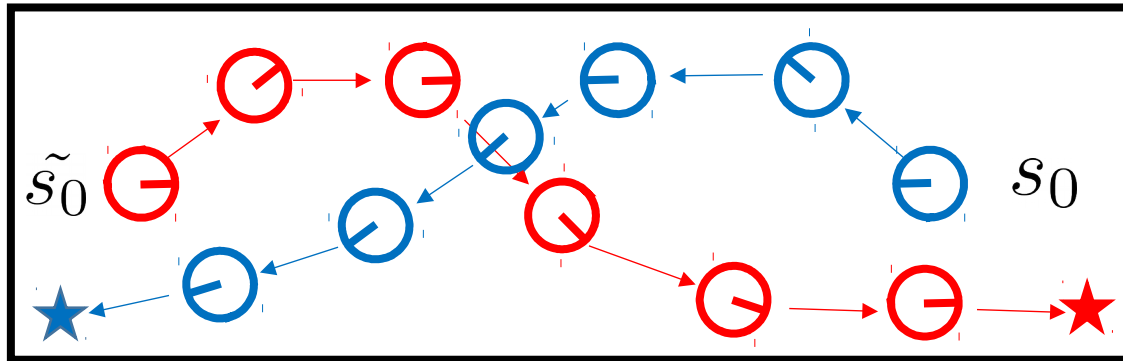


Algorithm 2: Deep V-learning

- 1 **Input:** trajectory training set D
- 2 **Output:** value network $V(\cdot; \mathbf{w})$
- 3 $V(\cdot; \mathbf{w}) \leftarrow \text{train_nn}(D)$ //step 1: initialization
- 4 duplicate value net $V' \leftarrow V$ //step 2: RL
- 5 initialize experience set $E \leftarrow D$
- 6 **for** $episode=1, \dots, N_{eps}$ **do**
- 7 **for** m times **do**
- 8 $s_0, \tilde{s}_0 \leftarrow \text{randomTestcase}()$
- 9 $s_{0:t_f} \leftarrow \text{CADRL}(V), \tilde{s}_{0:\tilde{t}_f} \leftarrow \text{CADRL}(V)$
- 10 $\mathbf{y}_{0:T}, \tilde{\mathbf{y}}_{0:\tilde{t}_f} \leftarrow \text{findValues}(V', s_{0:t_f}, \tilde{s}_{0:\tilde{t}_f})$
- 11 $E \leftarrow \text{assimilate}(E, (\mathbf{y}, \mathbf{s}^{jn})_{0:t_f}, (\tilde{\mathbf{y}}, \tilde{\mathbf{s}}^{jn})_{0:\tilde{t}_f})$
- 12 $e \leftarrow \text{randSubset}(E)$
- 13 $\mathbf{w} \leftarrow \text{backprop}(e)$
- 14 **for every** C episodes **do**
- 15 Evaluate(V), $V' \leftarrow V$
- 16 **return** V

Collision Avoidance Deep Reinforcement Learning

1. Train again with Deep reinforcement Learning



Algorithm 2: Deep V-learning

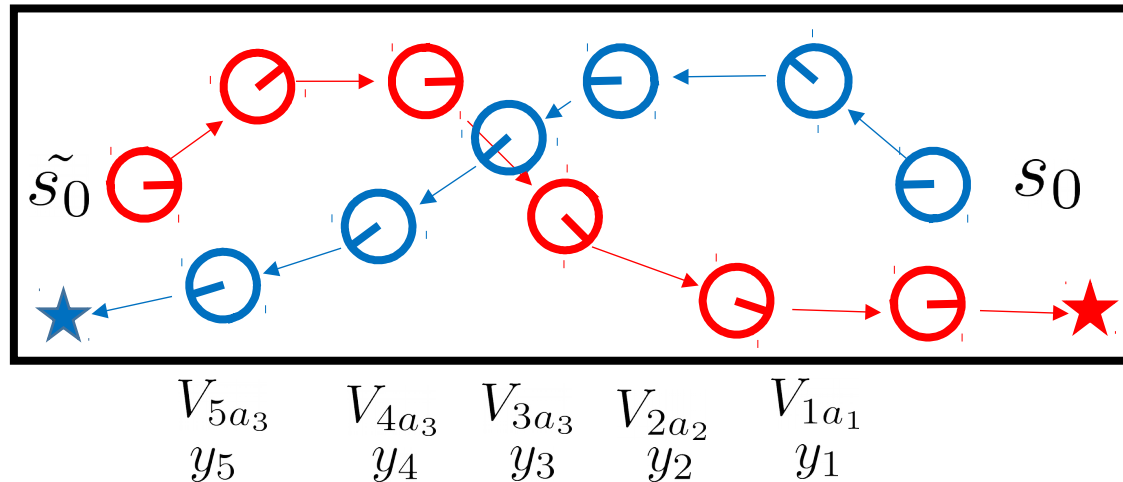
```

1 Input: trajectory training set  $D$ 
2 Output: value network  $V(\cdot; \mathbf{w})$ 
3  $V(\cdot; \mathbf{w}) \leftarrow \text{train\_nn}(D)$  //step 1: initialization
4 duplicate value net  $V' \leftarrow V$  //step 2: RL
5 initialize experience set  $E \leftarrow D$ 
6 for  $episode=1, \dots, N_{eps}$  do
7   for  $m$  times do
8      $s_0, \tilde{s}_0 \leftarrow \text{randomTestcase}()$ 
9      $s_{0:t_f} \leftarrow \text{CADRL}(V), \tilde{s}_{0:\tilde{t}_f} \leftarrow \text{CADRL}(V)$ 
10     $\mathbf{y}_{0:T}, \tilde{\mathbf{y}}_{0:\tilde{t}_f} \leftarrow \text{findValues}(V', s_{0:t_f}, \tilde{s}_{0:\tilde{t}_f})$ 
11     $E \leftarrow \text{assimilate}(E, (\mathbf{y}, \mathbf{s}^{jn})_{0:t_f}, (\tilde{\mathbf{y}}, \tilde{\mathbf{s}}^{jn})_{0:\tilde{t}_f})$ 
12   $e \leftarrow \text{randSubset}(E)$ 
13   $\mathbf{w} \leftarrow \text{backprop}(e)$ 
14  for every  $C$  episodes do
15    Evaluate( $V$ ),  $V' \leftarrow V$ 
16 return  $V$ 

```

Collision Avoidance Deep Reinforcement Learning

1. Train again with Deep reinforcement Learning



Algorithm 2: Deep V-learning

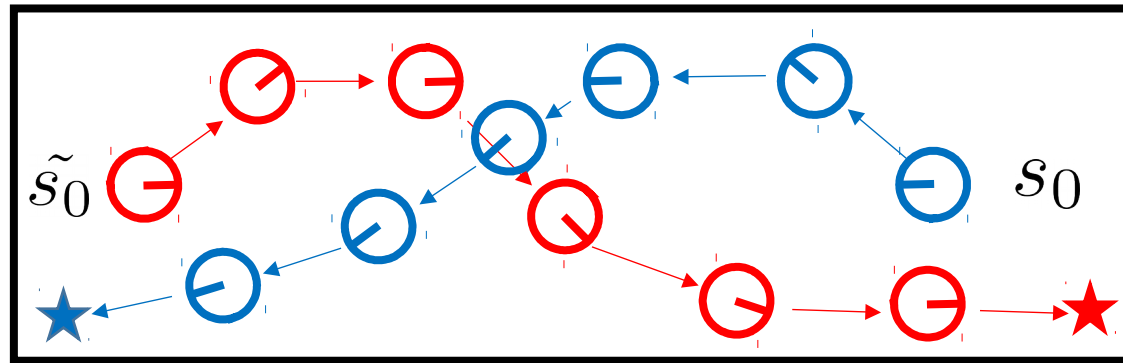
```

1 Input: trajectory training set  $D$ 
2 Output: value network  $V(\cdot; \mathbf{w})$ 
3  $V(\cdot; \mathbf{w}) \leftarrow \text{train\_nn}(D)$  //step 1: initialization
4 duplicate value net  $V' \leftarrow V$  //step 2: RL
5 initialize experience set  $E \leftarrow D$ 
6 for  $episode=1, \dots, N_{eps}$  do
7   for  $m$  times do
8      $s_0, \tilde{s}_0 \leftarrow \text{randomTestcase}()$ 
9      $s_{0:t_f} \leftarrow \text{CADRL}(V), \tilde{s}_{0:\tilde{t}_f} \leftarrow \text{CADRL}(V)$ 
10     $\mathbf{y}_{0:T}, \tilde{\mathbf{y}}_{0:\tilde{t}_f} \leftarrow \text{findValues}(V', s_{0:t_f}, \tilde{s}_{0:\tilde{t}_f})$ 
11     $E \leftarrow \text{assimilate}(E, (\mathbf{y}, \mathbf{s}^{jn})_{0:t_f}, (\tilde{\mathbf{y}}, \tilde{\mathbf{s}}^{jn})_{0:\tilde{t}_f})$ 
12     $e \leftarrow \text{randSubset}(E)$ 
13     $\mathbf{w} \leftarrow \text{backprop}(e)$ 
14    for every  $C$  episodes do
15      Evaluate( $V$ ),  $V' \leftarrow V$ 
16 return  $V$ 

```

Collision Avoidance Deep Reinforcement Learning

1. Train again with Deep reinforcement Learning



$$\begin{array}{ccccc} V_{5a_3} & V_{4a_3} & V_{3a_3} & V_{2a_2} & V_{1a_1} \\ y_5 & y_4 & y_3 & y_2 & y_1 \end{array}$$

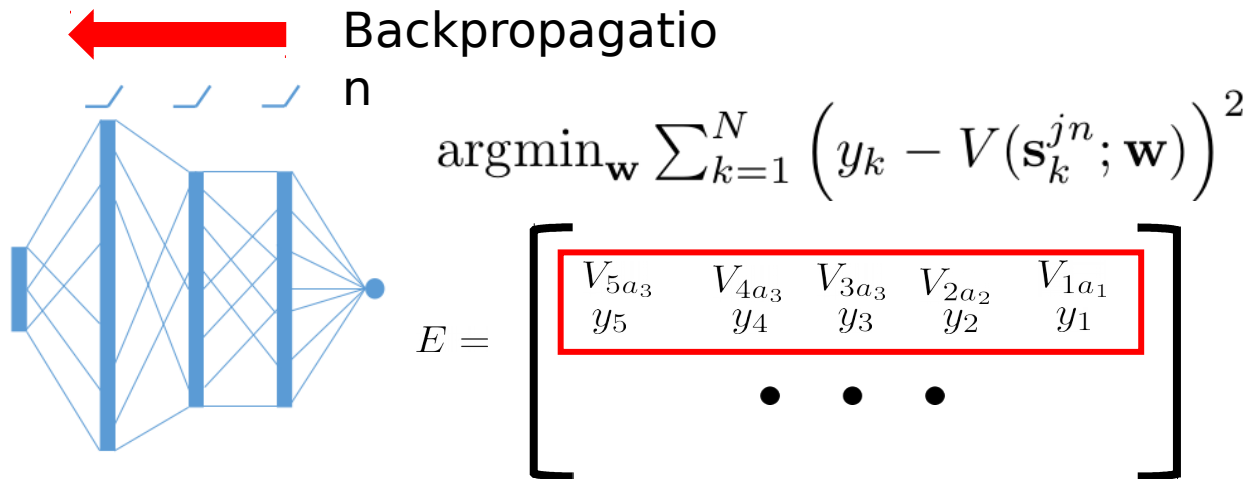
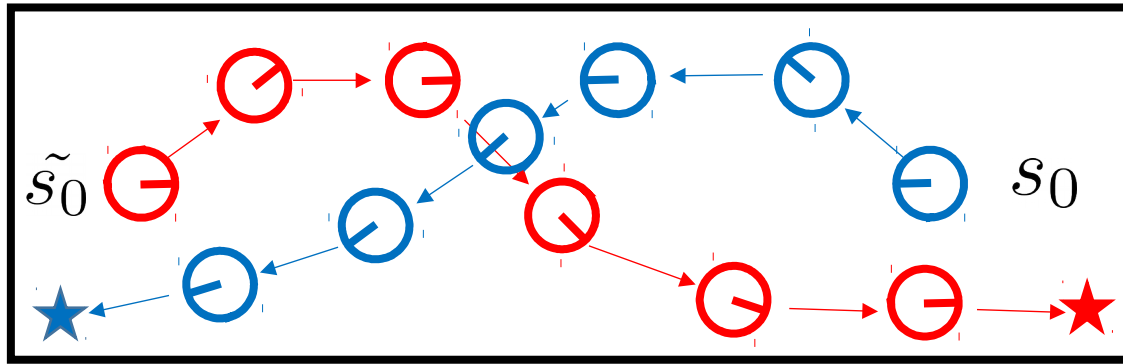
$$E = \left[\begin{array}{ccccc} V_{5a_3} & V_{4a_3} & V_{3a_3} & V_{2a_2} & V_{1a_1} \\ y_5 & y_4 & y_3 & y_2 & y_1 \\ \bullet & \bullet & \bullet & & \end{array} \right]$$

Algorithm 2: Deep V-learning

- 1 **Input:** trajectory training set D
- 2 **Output:** value network $V(\cdot; \mathbf{w})$
- 3 $V(\cdot; \mathbf{w}) \leftarrow \text{train_nn}(D)$ //step 1: initialization
- 4 duplicate value net $V' \leftarrow V$ //step 2: RL
- 5 initialize experience set $E \leftarrow D$
- 6 **for** $episode=1, \dots, N_{eps}$ **do**
- 7 **for** m times **do**
- 8 $s_0, \tilde{s}_0 \leftarrow \text{randomTestcase}()$
- 9 $s_{0:t_f} \leftarrow \text{CADRL}(V), \tilde{s}_{0:\tilde{t}_f} \leftarrow \text{CADRL}(V)$
- 10 $\mathbf{y}_{0:T}, \tilde{\mathbf{y}}_{0:\tilde{t}_f} \leftarrow \text{findValues}(V', s_{0:t_f}, \tilde{s}_{0:\tilde{t}_f})$
- 11 $E \leftarrow \text{assimilate}(E, (\mathbf{y}, \mathbf{s}^{jn})_{0:t_f}, (\tilde{\mathbf{y}}, \tilde{\mathbf{s}}^{jn})_{0:\tilde{t}_f})$
- 12 $e \leftarrow \text{randSubset}(E)$
- 13 $\mathbf{w} \leftarrow \text{backprop}(e)$
- 14 **for every** C episodes **do**
- 15 Evaluate(V), $V' \leftarrow V$
- 16 **return** V

Collision Avoidance Deep Reinforcement Learning

1. Train again with Deep reinforcement Learning

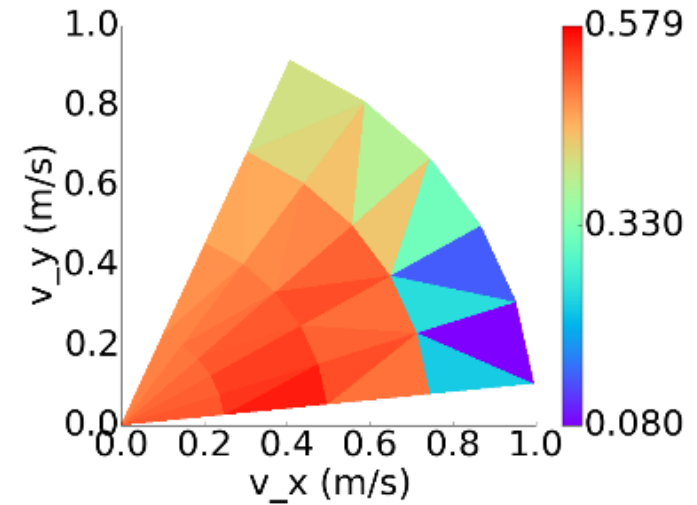
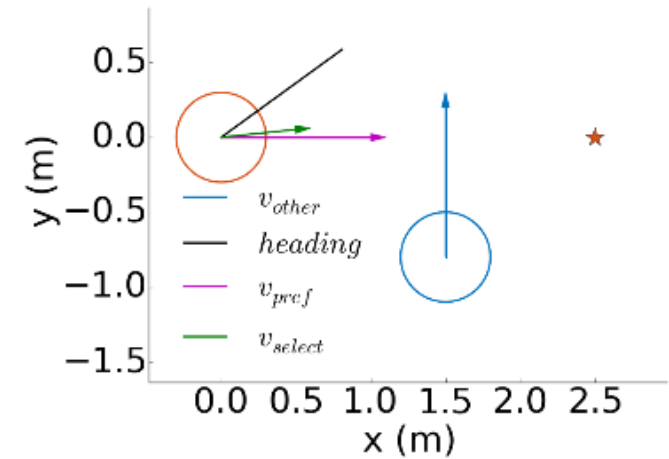
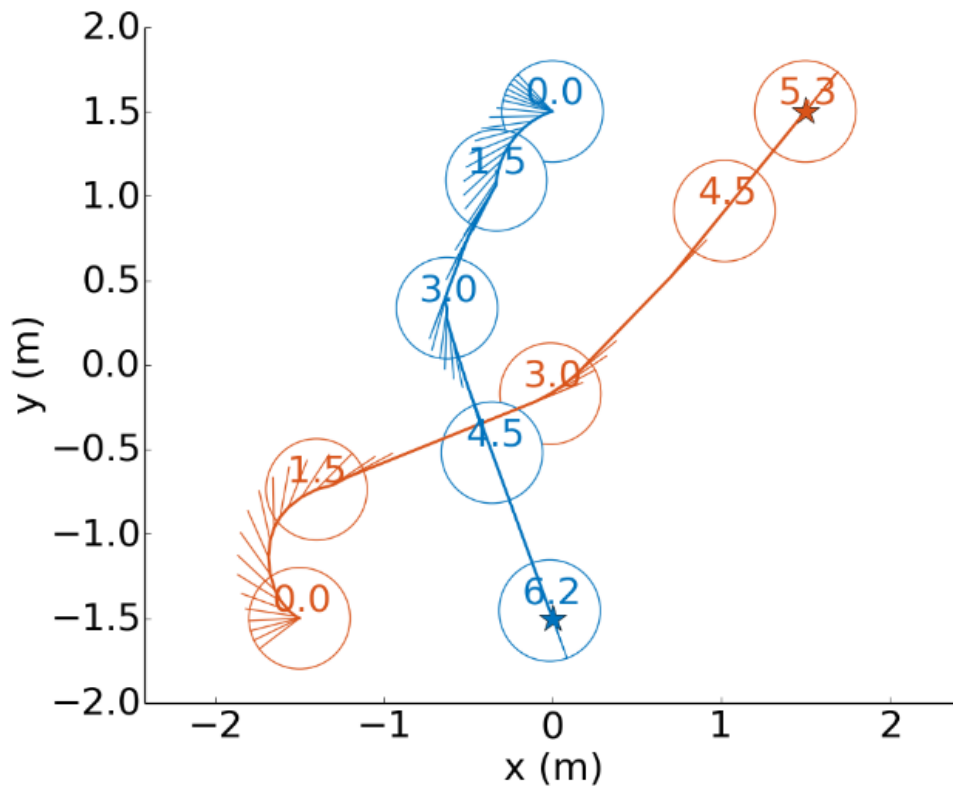


Algorithm 2: Deep V-learning

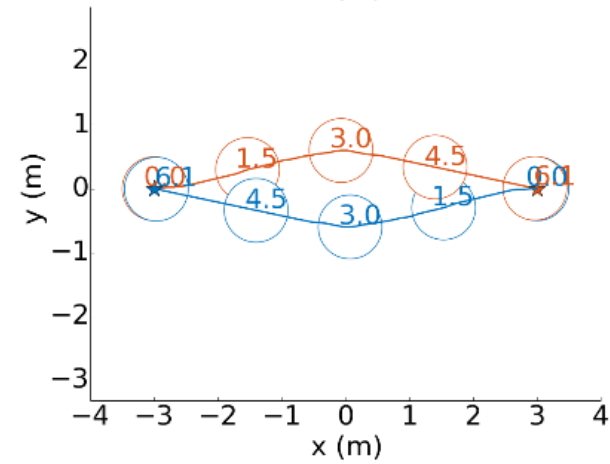
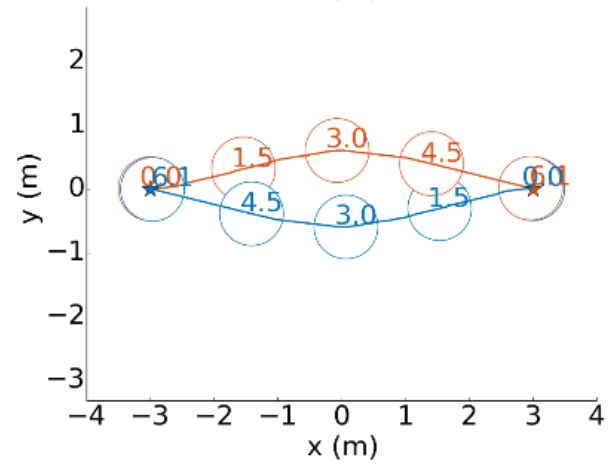
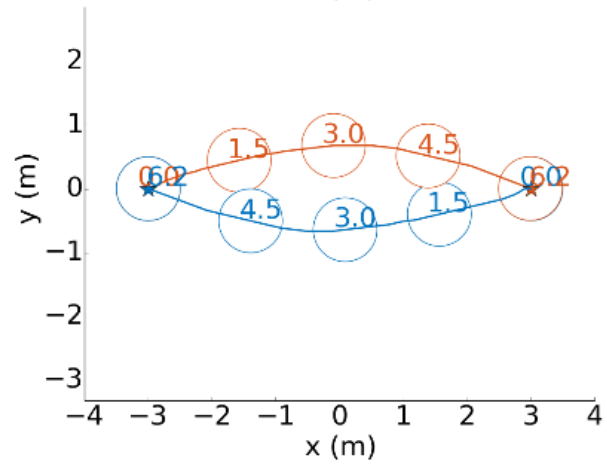
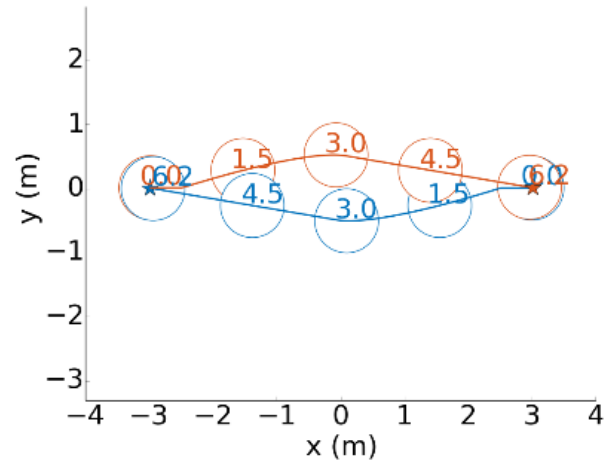
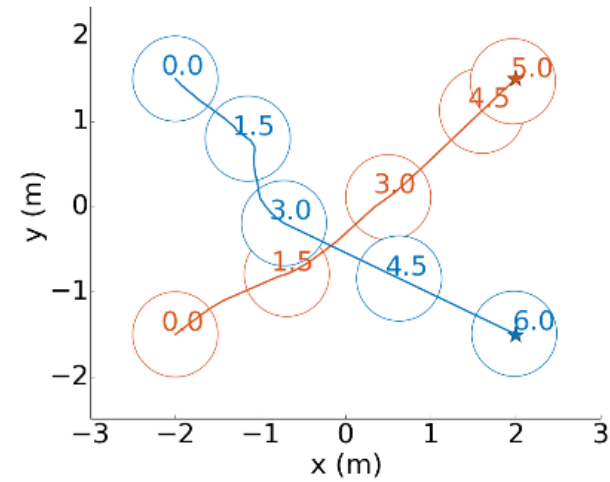
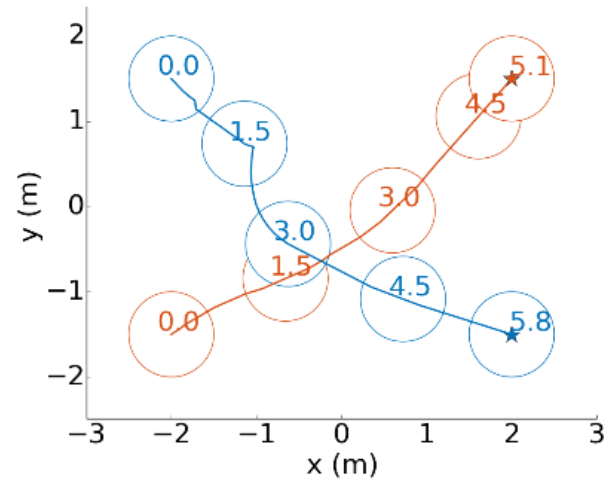
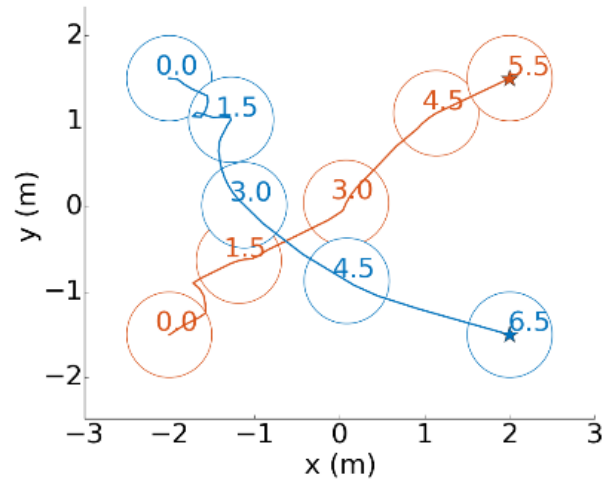
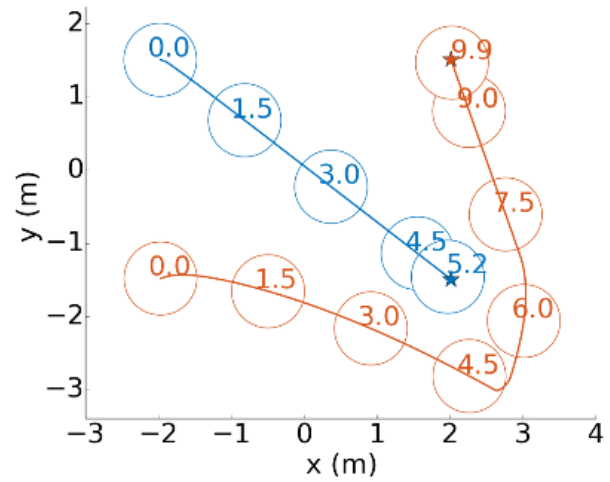
- 1 **Input:** trajectory training set D
- 2 **Output:** value network $V(\cdot; \mathbf{w})$
- 3 $V(\cdot; \mathbf{w}) \leftarrow \text{train_nn}(D)$ //step 1: initialization
- 4 duplicate value net $V' \leftarrow V$ //step 2: RL
- 5 initialize experience set $E \leftarrow D$
- 6 **for** $episode=1, \dots, N_{eps}$ **do**
- 7 **for** m times **do**
- 8 $\mathbf{s}_0, \tilde{\mathbf{s}}_0 \leftarrow \text{randomTestcase}()$
- 9 $\mathbf{s}_{0:t_f} \leftarrow \text{CADRL}(V), \tilde{\mathbf{s}}_{0:\tilde{t}_f} \leftarrow \text{CADRL}(V)$
- 10 $\mathbf{y}_{0:T}, \tilde{\mathbf{y}}_{0:\tilde{t}_f} \leftarrow \text{findValues}(V', \mathbf{s}_{0:t_f}, \tilde{\mathbf{s}}_{0:\tilde{t}_f})$
- 11 $E \leftarrow \text{assimilate}(E, (\mathbf{y}, \mathbf{s}^{jn})_{0:t_f}, (\tilde{\mathbf{y}}, \tilde{\mathbf{s}}^{jn})_{0:\tilde{t}_f})$
- 12 $e \leftarrow \text{randSubset}(E)$
- 13 $\mathbf{w} \leftarrow \text{backprop}(e)$
- 14 **for every** C episodes **do**
- 15 Evaluate(V), $V' \leftarrow V$
- 16 **return** V

Collision Avoidance Deep Reinforcement Learning

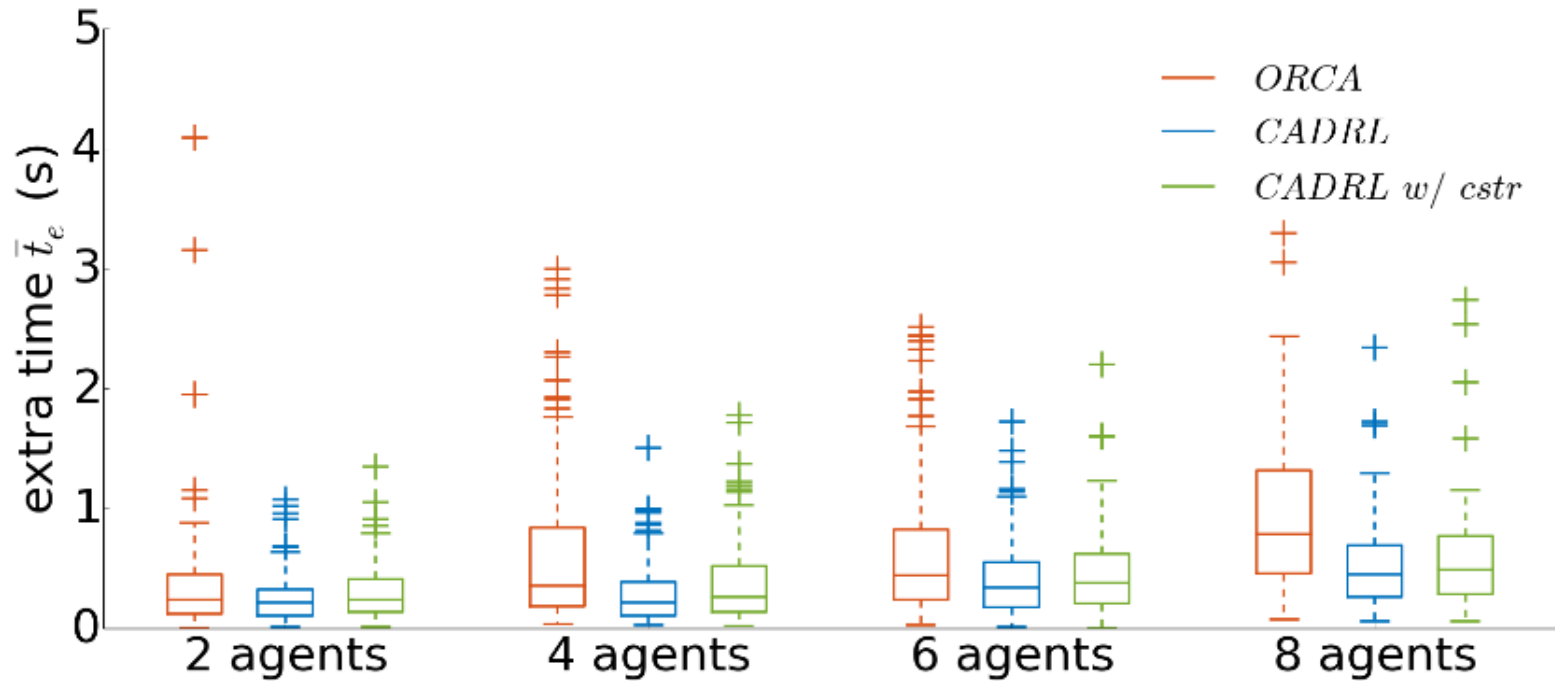
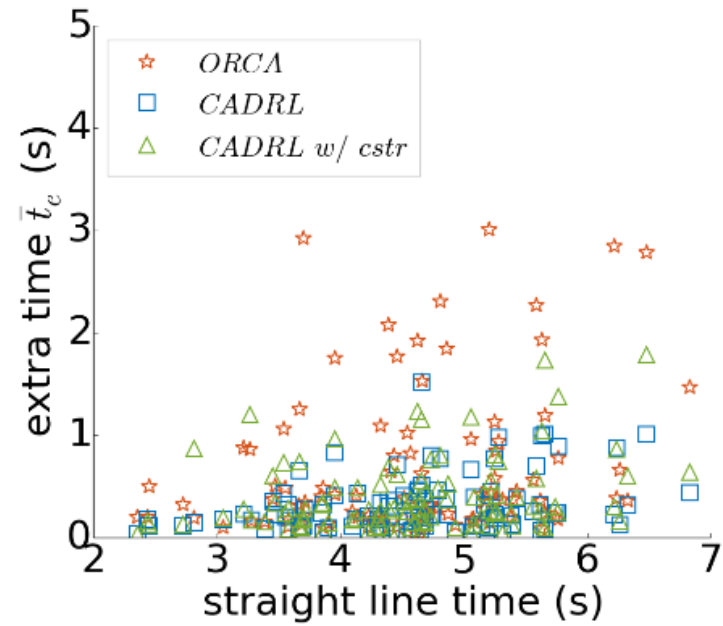
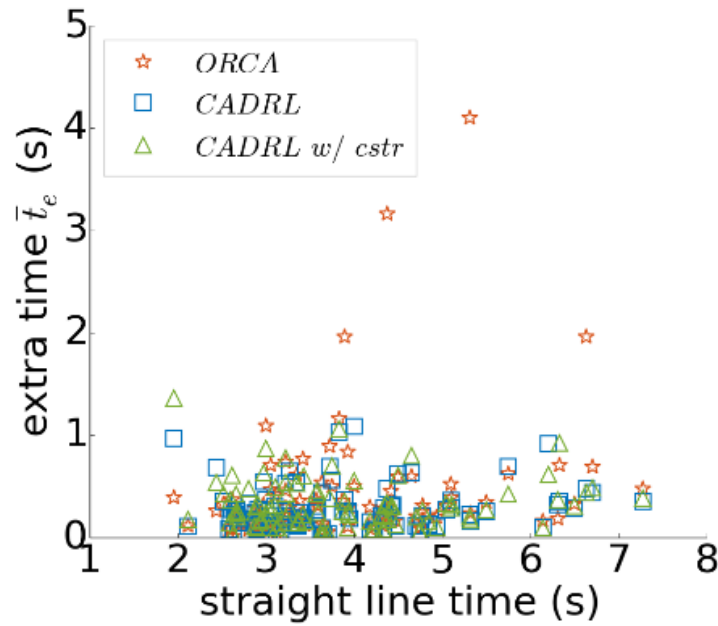
1. Train again with Deep reinforcement Learning



Result



Result



Result

Test case configuration		Extra time to goal \bar{t}_e (s) [Avg / 75th / 90th percentile]			Average min separation dist. (m)		
num agents	domain size (m)	ORCA	CADRL	CADRL w/ cstr	OCRA	CADRL	CADRL w/ cstr
2	4.0 × 4.0	0.46 / 0.45 / 0.73	0.27 / 0.33 / 0.56	0.31 / 0.42 / 0.60	0.122	0.199	0.198
4	5.0 × 5.0	0.69 / 0.85 / 1.85	0.31 / 0.40 / 0.76	0.39 / 0.53 / 0.86	0.120	0.192	0.191
6	6.0 × 6.0	0.65 / 0.83 / 1.50	0.44 / 0.56 / 0.87	0.48 / 0.63 / 1.02	0.118	0.117	0.180
8	7.0 × 7.0	0.96 / 1.33 / 1.84	0.54 / 0.70 / 1.01	0.59 / 0.77 / 1.09	0.110	0.171	0.170

Q&A

Quiz

- Values are update after each episode (T/F)
- Value function needs to be trained with ORCA (T/F)
- ORCA path does not need to be optimal (T/F)