

# DARA: Dynamics-Aware Reward Augmentation in Offline Reinforcement Learning

Presenter: Jinxin Liu (liujinxin@westlake.edu.cn)

Co-Authors: Hongyin Zhang, Donglin Wang

# Offline Reinforcement Learning (RL)

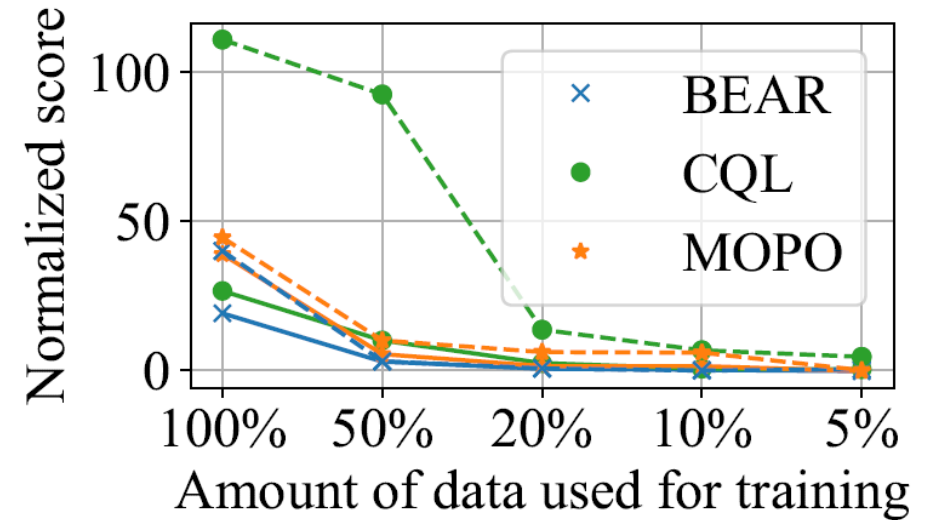
Offline reinforcement learning:  
learning from the previously collected dataset.

Offline-data-hungry!

Collecting a large offline dataset for one specific task over one specific environment is costly and laborious.

Offline Domain (Dynamics) Adaptation

- Source  $\rightarrow$  Target



# Offline Domain (Dynamics) Adaptation

Limited target offline data:  $\mathcal{D}$

$$\{(\mathbf{s}, \mathbf{a}, r, \mathbf{s}') \sim d_{\mathcal{D}}(\mathbf{s})\pi_b(\mathbf{a}|\mathbf{s})r(\mathbf{s}, \mathbf{a})T(\mathbf{s}'|\mathbf{s}, \mathbf{a})\}$$

Source offline data:  $\mathcal{D}'$

$$\{(\mathbf{s}, \mathbf{a}, r, \mathbf{s}') \sim d_{\mathcal{D}'}(\mathbf{s})\pi_{b'}(\mathbf{a}|\mathbf{s})r(\mathbf{s}, \mathbf{a})T'(\mathbf{s}'|\mathbf{s}, \mathbf{a})\}$$

Assumption:

- Same state space
- Same action space
- Same reward function
- Different transition dynamics
- Deterministic transition dynamics

**Definition 2** (Dynamics shift) Let  $\hat{M} := (\mathcal{S}, \mathcal{A}, r, \hat{T}, \rho_0, \gamma)$  be the empirical MDP estimated from  $\mathcal{D}$ . To evaluate a policy  $\pi$  for  $M := (\mathcal{S}, \mathcal{A}, r, T, \rho_0, \gamma)$  with offline dataset  $\mathcal{D}$ , we say that the dynamics shift (between  $\mathcal{D}$  and  $M$ ) in offline RL happens if there exists at least one transition pair  $(\mathbf{s}, \mathbf{a}, \mathbf{s}') \in \{(\mathbf{s}, \mathbf{a}, \mathbf{s}') : d_{\hat{M}}^{\pi}(\mathbf{s})\pi(\mathbf{a}|\mathbf{s})\hat{T}(\mathbf{s}'|\mathbf{s}, \mathbf{a}) > 0\}$  such that  $\hat{T}(\mathbf{s}'|\mathbf{s}, \mathbf{a}) \neq T(\mathbf{s}'|\mathbf{s}, \mathbf{a})$ .

**Lemma 2** Dynamics shift produces that  $\mathcal{B}_{\mathcal{D}}^{\pi}Q(\mathbf{s}, \mathbf{a}) \neq \mathcal{B}_{M}^{\pi}Q(\mathbf{s}, \mathbf{a})$  for some  $(\mathbf{s}, \mathbf{a})$  in  $\mathcal{S}'_{\pi}$ .

# Dynamics-Aware Reward Augmentation

**Lemma 2** *Dynamics shift produces that  $\mathcal{B}_{\mathcal{D}'}^{\pi} Q(\mathbf{s}, \mathbf{a}) \neq \mathcal{B}_M^{\pi} Q(\mathbf{s}, \mathbf{a})$  for some  $(\mathbf{s}, \mathbf{a})$  in  $S'_{\pi}$ .*

Resort an additional compensation  $\Delta_{\hat{T}', T}$  such that

$$\mathcal{B}_{\mathcal{D}'}^{\pi} Q(\mathbf{s}, \mathbf{a}) + \Delta_{\hat{T}', T}(\mathbf{s}, \mathbf{a}) = \mathcal{B}_M^{\pi} Q(\mathbf{s}, \mathbf{a})$$

---

## Algorithm 1 Framework for Dynamics-Aware Reward Augmentation (DARA)

---

**Require:** Target offline data  $\mathcal{D}$  (reduced) and source offline data  $\mathcal{D}'$

- 1: Learn classifiers ( $q_{\text{sas}}$  and  $q_{\text{sa}}$ ) that distinguish source data  $\mathcal{D}'$  from target data  $\mathcal{D}$ . (See Appendix [A.1.3](#))
  - 2: Set dynamics-aware  $\Delta r(\mathbf{s}_t, \mathbf{a}_t, \mathbf{s}_{t+1}) = \log \frac{q_{\text{sas}}(\text{source}|\mathbf{s}_t, \mathbf{a}_t, \mathbf{s}_{t+1})}{q_{\text{sas}}(\text{target}|\mathbf{s}_t, \mathbf{a}_t, \mathbf{s}_{t+1})} - \log \frac{q_{\text{sa}}(\text{source}|\mathbf{s}_t, \mathbf{a}_t)}{q_{\text{sa}}(\text{target}|\mathbf{s}_t, \mathbf{a}_t)}$ .
  - 3: Modify rewards for all  $(\mathbf{s}_t, \mathbf{a}_t, r_t, \mathbf{s}_{t+1})$  in  $\mathcal{D}'$ :  $r_t \leftarrow r_t - \eta \Delta r$ .
  - 4: Learn policy with  $\{\mathcal{D} \cup \mathcal{D}'\}$  using prior model-free or model-based offline RL algorithms.
-

# Dynamics-Aware Reward Augmentation

**Lemma 2** *Dynamics shift produces that  $\mathcal{B}_{\mathcal{D}'}^\pi Q(\mathbf{s}, \mathbf{a}) \neq \mathcal{B}_M^\pi Q(\mathbf{s}, \mathbf{a})$  for some  $(\mathbf{s}, \mathbf{a})$  in  $S'_\pi$ .*

Resort an additional compensation  $\Delta_{\hat{T}', T}$  such that

$$\mathcal{B}_{\mathcal{D}'}^\pi Q(\mathbf{s}, \mathbf{a}) + \Delta_{\hat{T}', T}(\mathbf{s}, \mathbf{a}) = \mathcal{B}_M^\pi Q(\mathbf{s}, \mathbf{a})$$

---

## Algorithm 1 Framework for Dynamics-Aware Reward Augmentation (DARA)

---

**Require:** Target offline data  $\mathcal{D}$  (reduced) and source offline data  $\mathcal{D}'$

- 1: Learn classifiers ( $q_{\text{sas}}$  and  $q_{\text{sa}}$ ) that distinguish source data  $\mathcal{D}'$  from target data  $\mathcal{D}$ . (See Appendix [A.1.3](#))
  - 2: Set dynamics-aware  $\Delta r(\mathbf{s}_t, \mathbf{a}_t, \mathbf{s}_{t+1}) = \log \frac{q_{\text{sas}}(\text{source}|\mathbf{s}_t, \mathbf{a}_t, \mathbf{s}_{t+1})}{q_{\text{sas}}(\text{target}|\mathbf{s}_t, \mathbf{a}_t, \mathbf{s}_{t+1})} - \log \frac{q_{\text{sa}}(\text{source}|\mathbf{s}_t, \mathbf{a}_t)}{q_{\text{sa}}(\text{target}|\mathbf{s}_t, \mathbf{a}_t)}$ .
  - 3: Modify rewards for all  $(\mathbf{s}_t, \mathbf{a}_t, r_t, \mathbf{s}_{t+1})$  in  $\mathcal{D}'$ :  $r_t \leftarrow r_t - \eta \Delta r$ .
  - 4: Learn policy with  $\{\mathcal{D} \cup \mathcal{D}'\}$  using prior model-free or model-based offline RL algorithms.
- 

$\Delta_{\hat{T}', T}$  discourages the learning from these offline transitions that are likely in source but are unlikely in target.

# Experiments

Body Mass Shift		10T	1T	1T+10S w/o Aug.	1T+10S DARA	10T	1T	1T+10S w/o Aug.	1T+10S DARA	10T	1T	1T+10S w/o Aug.	1T+10S DARA
Hopper		BEAR				BRAC-p				AWR			
	Random	11.4	1.0 ↓	4.6 ↑	8.4 ↑	11.0	10.9 ↓	9.6 ↓	<b>11.0</b> ↑	10.2	10.3 ↑	3.4 ↓	4.5 ↑
	Medium	52.1	0.8 ↓	0.9 ↑	1.6 ↑	32.7	29.0 ↓	29.2 ↑	<b>32.9</b> ↑	35.9	30.9 ↓	20.8 ↓	28.9 ↑
	Medium-R	33.7	1.3 ↓	18.2 ↑	<b>34.1</b> ↑	0.6	5.4 ↑	20.1 ↑	<b>30.8</b> ↑	28.4	8.8 ↓	4.1 ↓	4.2 ↑
Medium-E	96.3	0.8 ↓	0.6 ↓	1.2 ↑	1.9	34.5 ↑	32.3 ↓	<b>34.7</b> ↑	27.1	27.0 ↓	26.8 ↓	<b>26.6</b> ↓	
Hopper		BCQ				CQL				MOPO			
	Random	10.6	10.6 ↓	8.3 ↓	<b>9.7</b> ↑	10.8	10.6 ↓	10.2 ↓	<b>10.4</b> ↑	11.7	4.8 ↓	2.0 ↓	2.1 ↑
	Medium	54.5	37.1 ↓	25.7 ↓	38.4 ↑	58.0	43.0 ↓	44.9 ↑	<b>59.3</b> ↑	28.0	4.1 ↓	5.0 ↑	10.7 ↑
	Medium-R	33.1	9.3 ↓	28.7 ↑	<b>32.8</b> ↑	48.6	9.6 ↓	1.4 ↓	3.7 ↑	67.5	1.0 ↓	5.5 ↑	8.4 ↑
Medium-E	110.9	58 ↓	75.4 ↑	84.2 ↑	98.7	59.7 ↓	53.6 ↓	<b>99.7</b> ↑	23.7	1.6 ↓	4.8 ↑	5.8 ↑	
Walker2d		BEAR				BRAC-p				AWR			
	Random	7.3	1.5 ↓	3.1 ↑	3.2 ↑	-0.2	0.0 ↑	1.3 ↑	<b>3.2</b> ↑	1.5	1.3 ↓	2.0 ↑	<b>2.4</b> ↑
	Medium	59.1	-0.5 ↓	0.6 ↑	0.3 ↓	77.5	6.4 ↓	70.0 ↑	<b>78.0</b> ↑	17.4	14.8 ↓	17.1 ↑	<b>17.2</b> ↑
	Medium-R	19.2	0.7 ↓	6.5 ↑	7.3 ↑	-0.3	8.5 ↑	9.9 ↑	<b>18.6</b> ↑	15.5	7.4 ↓	1.6 ↓	1.5 ↓
Medium-E	40.1	-0.1 ↓	1.5 ↑	2.3 ↑	76.9	20.6 ↓	64.1 ↑	<b>77.5</b> ↑	53.8	35.5 ↓	52.5 ↑	<b>53.3</b> ↑	
Walker2d		BCQ				CQL				MOPO			
	Random	4.9	1.8 ↓	4.5 ↑	<b>4.8</b> ↑	7.0	1.7 ↓	3.2 ↑	3.4 ↑	13.6	-0.2 ↓	-0.1 ↑	-0.1 ↓
	Medium	53.1	32.8 ↓	50.9 ↑	<b>52.3</b> ↑	79.2	42.9 ↓	80.0 ↑	<b>81.7</b> ↑	17.8	7.0 ↓	5.7 ↓	11.0 ↑
	Medium-R	15.0	6.9 ↓	14.9 ↑	<b>15.1</b> ↑	26.7	4.6 ↓	0.8 ↓	2.0 ↑	39.0	5.1 ↓	3.1 ↓	14.2 ↑
Medium-E	57.5	32.5 ↓	55.2 ↑	<b>57.2</b> ↑	111.0	49.5 ↓	63.5 ↑	93.3 ↑	44.6	5.3 ↓	5.5 ↑	17.2 ↑	

DARA can enable an adaptive policy with reduced offline data in target.

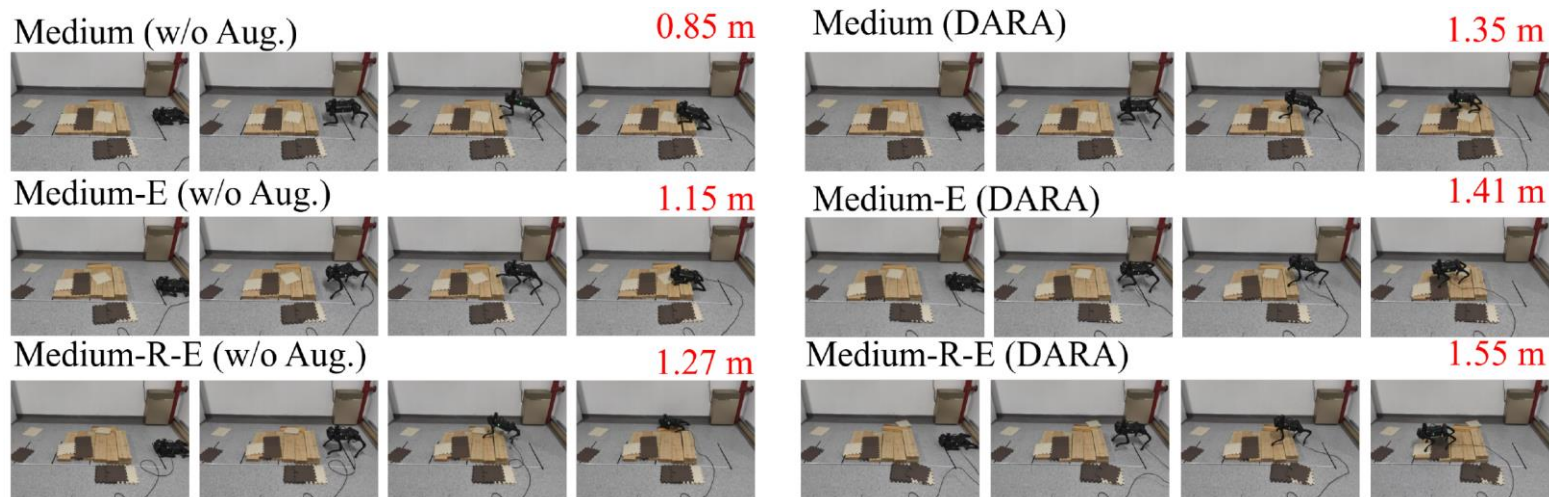
# Experiments

Body Mass Shift		Tune	DARA	Tune	DARA	Tune	DARA	Tune	DARA	Tune	DARA	$\pi_p \hat{T}$	$\hat{T} \pi_p$
Hopper		BEAR		BRAC-p		BCQ		CQL		MOPO		MABE	
	Random	0.8	8.4 $\uparrow$	6.0	11.0 $\uparrow$	8.8	9.7 $\uparrow$	<b>31.6</b>	10.4 $\downarrow$	0.7	2.1 $\uparrow$	10.6	9.0
	Medium	0.8	1.6 $\uparrow$	22.7	32.9 $\uparrow$	31.7	38.4 $\uparrow$	44.5	<b>59.3</b> $\uparrow$	0.7	10.7 $\uparrow$	48.8	23.1
	Medium-R	0.7	<b>34.1</b> $\uparrow$	14.7	30.8 $\uparrow$	27.5	32.8 $\uparrow$	1.3	3.7 $\uparrow$	0.6	8.4 $\uparrow$	17.1	20.4
	Medium-E	0.9	1.2 $\uparrow$	19.2	34.7 $\uparrow$	85.9	84.2 $\downarrow$	47.6	<b>99.7</b> $\uparrow$	2.2	5.8 $\uparrow$	28.1	38.9
Walker2d		BEAR		BRAC-p		BCQ		CQL		MOPO		MABE	
	Random	<b>6.6</b>	3.2 $\downarrow$	3.9	3.2 $\downarrow$	4.7	4.8 $\uparrow$	1.1	3.4 $\uparrow$	0.1	-0.1 $\downarrow$	6.0	-0.2
	Medium	0.3	0.3 $\downarrow$	76.0	78.0 $\uparrow$	28.4	52.3 $\uparrow$	72.3	<b>81.7</b> $\uparrow$	-0.2	11.0 $\uparrow$	30.1	56.7
	Medium-R	1.2	7.3 $\uparrow$	10.0	<b>18.6</b> $\uparrow$	10.4	15.1 $\uparrow$	1.8	2.0 $\uparrow$	0.0	14.2 $\uparrow$	13.3	12.5
	Medium-E	2.4	2.3 $\downarrow$	74.5	77.5 $\uparrow$	22.7	57.2 $\uparrow$	68.6	<b>93.3</b> $\uparrow$	7.3	17.2 $\uparrow$	43.7	82.7

Comparison with cross-domain baselines.

(BCQ)	w/o Aug.	DARA
Medium	0.85	1.35 $\uparrow$
Medium-E	1.15	1.41 $\uparrow$
Medium-R-E	1.27	1.55 $\uparrow$

Sim2real: Deployment on the obstructive and dynamic environment of BCQ.





## Conclusion

1. The characterization of the dynamics shift in offline RL and the derivation of dynamics-aware reward augmentation (DARA) framework built on prior model-free and model-based formulations.
2. With only modest amounts of target offline data, we show that DARA-based offline methods can acquire an adaptive policy for the target tasks and achieve better performance compared to baselines in both simulated and real-world tasks.