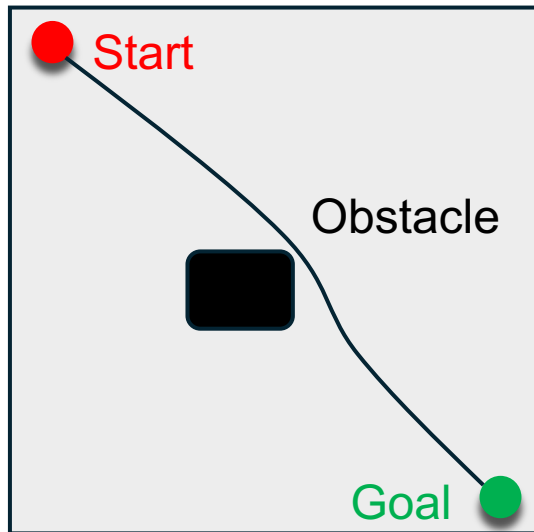


# Off-Dynamics Reinforcement Learning via Domain Adaptation and Reward Augmented Imitation

Yihong Guo, Yixuan Wang, Yuanyuan Shi, Pan Xu, Anqi Liu

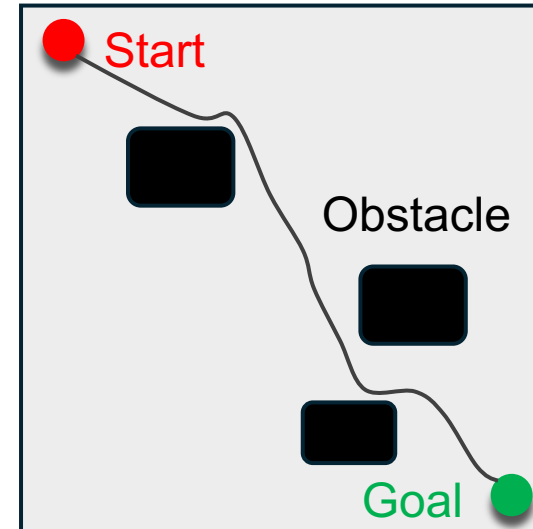


# Off-dynamics Reinforcement Learning



Source domain (simulation)

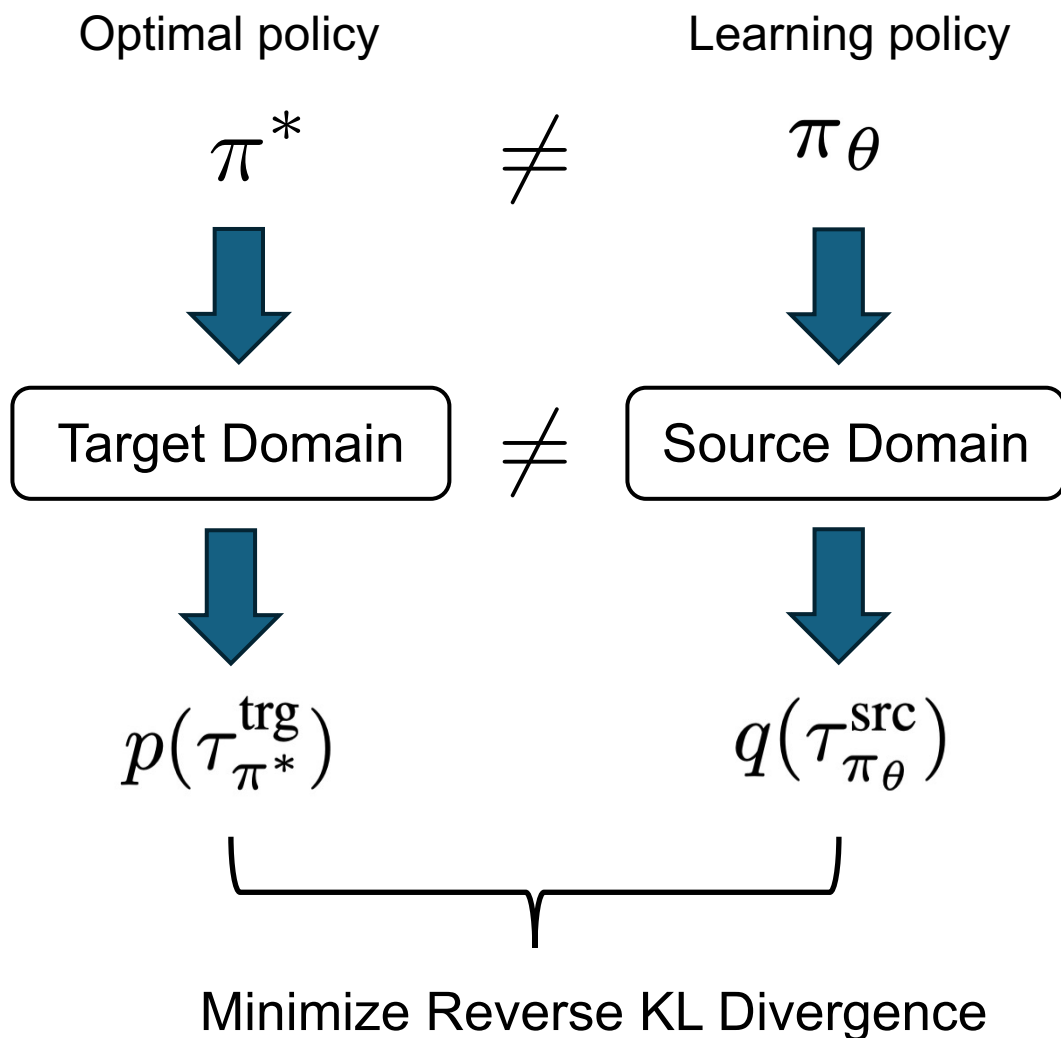
Train on source  
→  
Deploy to target



Target domain (real-world)

- Different transition probability:  $P_{src}(s_{t+1}|s_t, a_t) \neq P_{trg}(s_{t+1}|s_t, a_t)$
- Same state, action space and reward function.

# Limitation of existing reward shaping methods DARC [1]

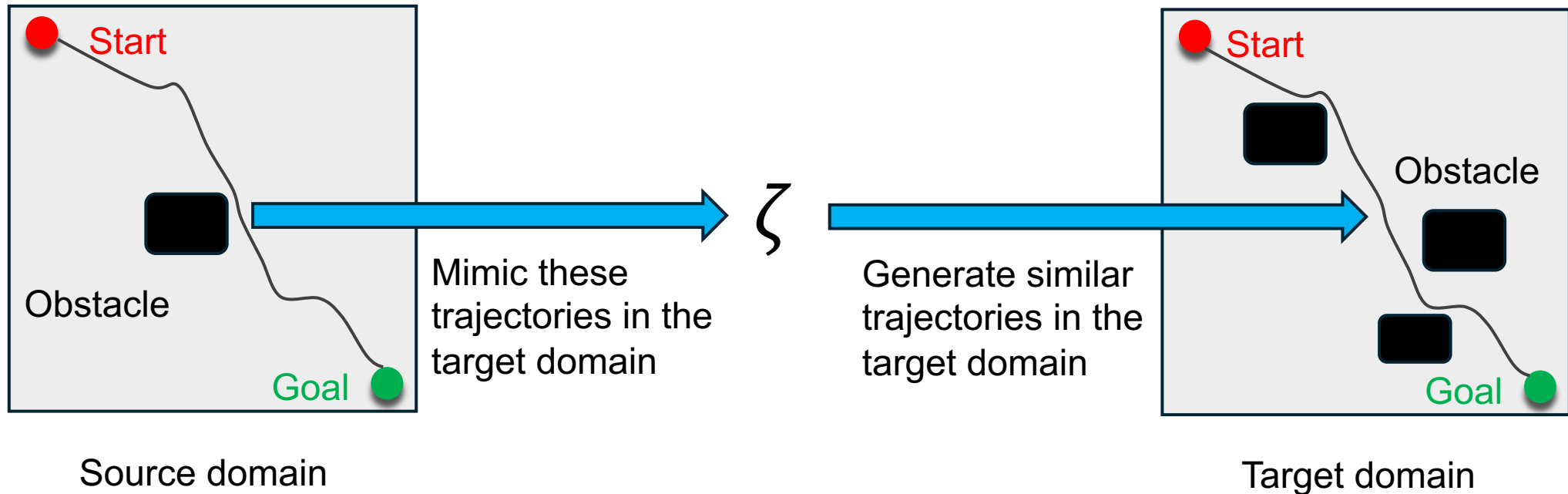


- DARC policy is suboptimal in the target domain.
- DARC relies on the assumption of target optimal policy performs well in the source domain.

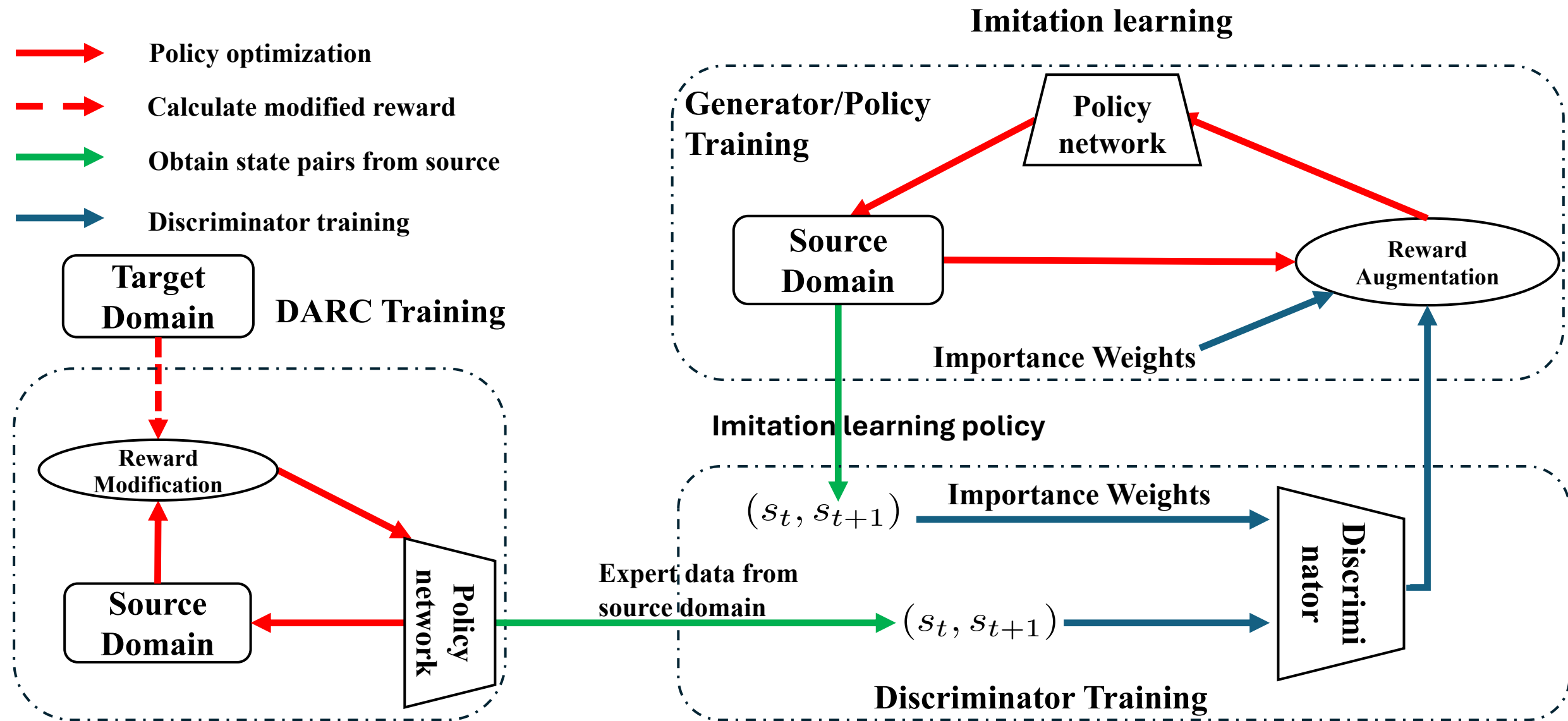




# Domain Adaptation and Reward Augmented Imitation Learning (DARAIL)



# Domain Adaptation and Reward Augmented Imitation Learning (DARAIL)



# Experiments

- Experiment on Mujoco: HalfCheetah, Ant, Walker2d and Reacher
- Dynamics shift:
  - Broken environment: set the 0-th index action to 0.
  - Modifying parameters: set the scale of gravity/density from 1.0 to 0.5/1.5 in the target domain.
- Evaluate with mean reward and standard deviation among multiple runs.

# Comparison with DARC

- DARC evaluation worse than DARC training.
- DARAIL outperforms DARC evaluation results.

Table 1: Comparison of DARAIL with DARC, broken source environment.

	DARC Evaluation	DARC Training	Optimal in Target	DARAIL
HalfCheetah	4133 $\pm$ 828	6995 $\pm$ 30	8543 $\pm$ 230	7067 $\pm$ 176
Ant	4280 $\pm$ 33	5197 $\pm$ 155	6183 $\pm$ 348	5357 $\pm$ 79
Walker2d	2669 $\pm$ 788	3896 $\pm$ 523	3899 $\pm$ 214	4366 $\pm$ 434
Reacher	-26.3 $\pm$ 3.3	-11.2 $\pm$ 2.9	-7.2 $\pm$ 1.2	-13.7 $\pm$ 0.9

Table 2: Comparison of DARAIL with DARC, 1.5 gravity.

	DARC Evaluation	DARC Training	Optimal in Target	DARAIL
HalfCheetah	653 $\pm$ 142	4897 $\pm$ 653	6894 $\pm$ 491	4093 $\pm$ 1021
Ant	1587 $\pm$ 594	2170 $\pm$ 258	5320 $\pm$ 429	3472 $\pm$ 771
Walker2d	257 $\pm$ 28	4130 $\pm$ 689	4254 $\pm$ 345	4409 $\pm$ 401
Reacher	-55.3 $\pm$ 10.3	-17.2 $\pm$ 3.8	-8.3 $\pm$ 1.3	-9.5 $\pm$ 0.22



# DARAIL outperforms other baselines

Table 3: Comparison of DARAIL with baselines in off-dynamics RL, broken source environment.

	DAIL	IS-R	IS-ACL	MBPO	MATL	GARAT	DARAIL
HalfCheetah	$6402 \pm 362$	$6007 \pm 863$	$6934 \pm 231$	$4323 \pm 7$	$1538 \pm 616$	$5877 \pm 382$	<b><math>7067 \pm 176</math></b>
Ant	$3239 \pm 395$	$1463 \pm 1055$	$2753 \pm 94$	$2445 \pm 13$	$2006 \pm 17$	$3380 \pm 268$	<b><math>5357 \pm 79</math></b>
Walker2d	$2330 \pm 156$	$3092 \pm 434$	$3881 \pm 269$	$1012 \pm 41$	$250 \pm 5$	$3296 \pm 284$	<b><math>4366 \pm 434</math></b>
Reacher	$-13.9 \pm 1.1$	$-17.6 \pm 0.25$	$-14.1 \pm 0.16$	$-14.3 \pm 2$	$-30 \pm 10$	$-14.7 \pm 2.6$	<b><math>-13.7 \pm 0.9</math></b>

Table 4: Comparison of DARAIL with baselines in off-dynamics RL, 1.5 gravity.

	DAIL	IS-R	IS-ACL	MBPO	MATL	GARAT	DARAIL
HalfCheetah	$2666 \pm 2037$	$2718 \pm 1978$	$3576 \pm 312$	$619 \pm 311$	$337 \pm 205$	$3825 \pm 437$	<b><math>4093 \pm 1021</math></b>
Ant	$990 \pm 251$	$1712 \pm 393$	$2396 \pm 573$	$989 \pm 13$	$1376 \pm 466$	$1961 \pm 115$	<b><math>3472 \pm 771</math></b>
Walker2d	$525 \pm 142$	$1543 \pm 604$	$1369 \pm 705$	$870 \pm 451$	$1419 \pm 489$	$630 \pm 230$	<b><math>4409 \pm 401</math></b>
Reacher	$-16.5 \pm 1.1$	$-14.6 \pm 0.8$	$-47.4 \pm 8.3$	$-18.3 \pm 0.9$	$-17.6 \pm 0.7$	$-16.7 \pm 0.3$	<b><math>-9.5 \pm 0.22</math></b>

# Summary

- Propose DARAIL for off-dynamics RL.
- Recognize the limitations of DARC and other works with same reward shaping method.
- Propose imitation learning with augmented reward estimator to address the limitation of DARC.
- Propose an error bound with a relaxed assumption.