

# Reining Generalization in Offline Reinforcement Learning via Representation Distinction

---

**Yi Ma<sup>1</sup>, Hongyao Tang<sup>2,\*</sup>, Dong Li<sup>3</sup>, Zhaopeng Meng<sup>1</sup>**

1 College of Intelligence and Computing, Tianjin University

2 Université de Montréal, Mila

3 Noah's Ark Lab, Huawei Technology

\* Corresponding authors

# Outline

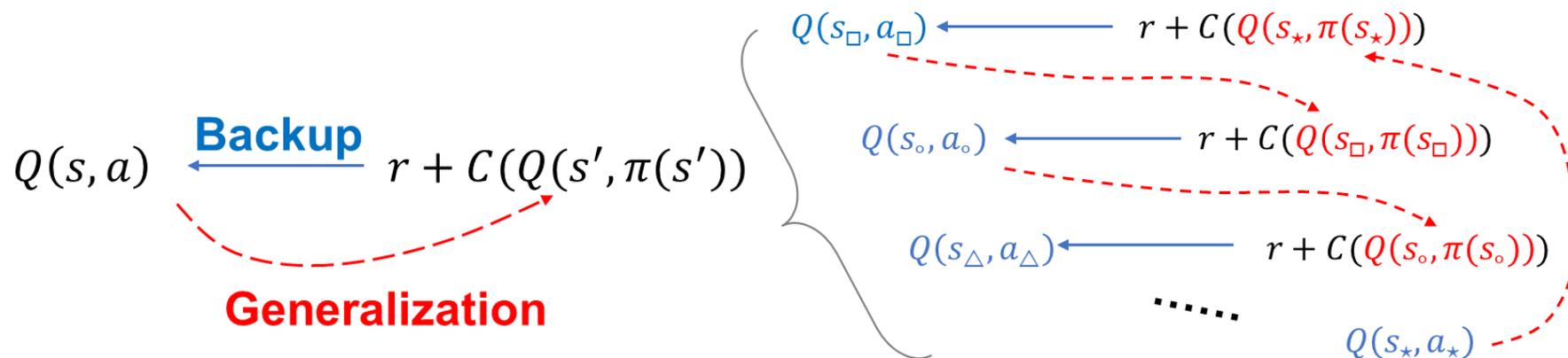
---



- **Problem Formulation**
- **Method**
- **Experimental Results**

## Problem Formulation: Backup-Generalization Cycle in Offline RL

- We introduce a view called Backup-Generalization Cycle. This view, as depicted in the Fig below, fosters an understanding of typical offline value function learning via two key components: **Backup** and **Generalization**:



This dynamic interplay forms a cycle:

- (1) the backups on  $(s, a) \in D$  consistently influence the values of  $(s, a) \notin D$  through generalization;
- (2) the consistently changing  $Q(s', \pi(s'))$  participates in the backups on  $(s, a) \notin D$ ;

The two kinds of dynamics iterate and twine during the learning process.

## Problem Formulation: Overgeneralization in Offline RL

---

Further, we consider how Q function update caused by typical Temporal-Difference (TD) learning on a single state-action pair  $(s, a) \in D$  (denoted as  $\phi \rightarrow \phi'$ ), affects the Q-value of an arbitrary state-action pair  $(\bar{s}, \bar{a})$ .

The post-update parameter  $\phi'$  can be formalized as follows:

$$\phi' = \phi + (\mathcal{T}Q_\phi(s, a) - Q_\phi(s, a))\nabla_\phi Q_\phi(s, a)$$

By Taylor expansion at the pre-update parameter  $\phi$ :

$$Q_{\phi'}(\bar{s}, \bar{a}) = Q_\phi(\bar{s}, \bar{a}) + \nabla_\phi Q_\phi(\bar{s}, \bar{a})^T (\phi' - \phi) + \mathcal{O}(\|\phi' - \phi\|^2)$$

By plugging the first Eq to the second Eq:

$$Q_{\phi'}(\bar{s}, \bar{a}) = Q_\phi(\bar{s}, \bar{a}) + k_\phi(\bar{s}, \bar{a}, s, a)(\mathcal{T}Q_\phi(s, a) - Q_\phi(s, a)) + \mathcal{O}(\|\phi' - \phi\|^2)$$

where  $k_\phi(\bar{s}, \bar{a}, s, a) = \nabla_\phi Q_\phi(\bar{s}, \bar{a})^T \nabla_\phi Q_\phi(s, a)$ , which we called Neural Tangent Kernel. We can control the generalization by mainly adjusting this kernel.

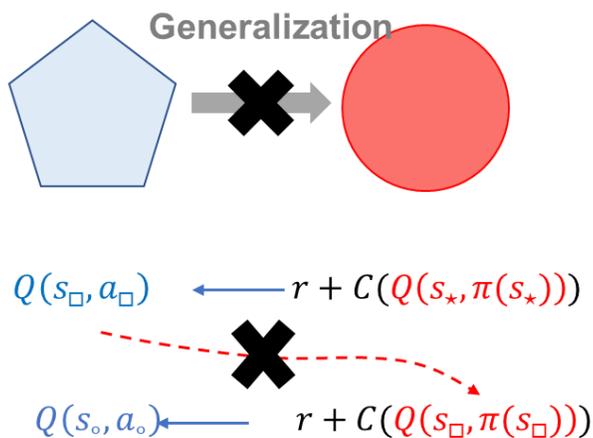
# Method: Reining Generalization via Two-Stage Kernel Control

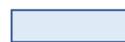
Suppress  $Q(s, a)$  generalization  $\rightarrow$   
 $Q(s, \pi(s))$  for  $(s, a) \in D$  by  
 $\min_{\phi} |\nabla_{\phi} Q_{\phi}(s, a)^T \nabla_{\phi} Q_{\phi}(s, \pi(s))|$

When the learning policy evolves and resembles the behavioral policy  $\min_{\phi} |\nabla_{\phi} Q_{\phi}(s, a)^T \nabla_{\phi} Q_{\phi}(s, \pi(s))|$  for  $a \approx \pi(s)$  could lead to over-inhibition

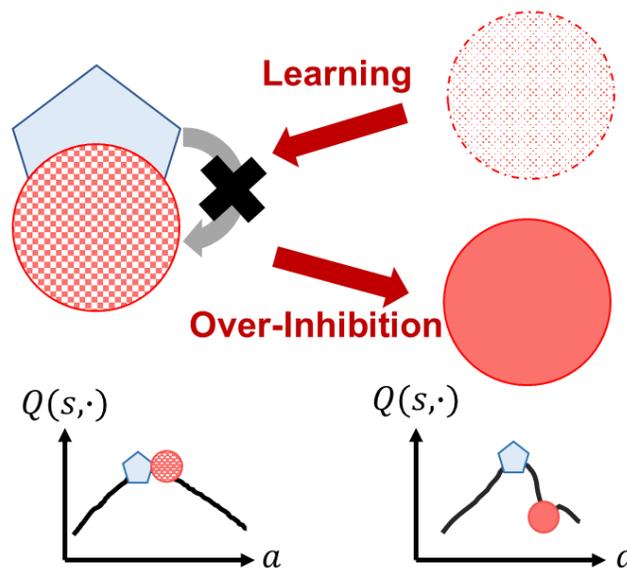
Suppress the generalization between the learning policy distribution and the designed OOD policy distribution by  $\min_{\phi} |\nabla_{\phi} Q_{\phi}(s, \pi(s))^T \nabla_{\phi} Q_{\phi}(s, \pi_{ood}(s))|$

**Policy-Dataset  
Generalization Inhibition**



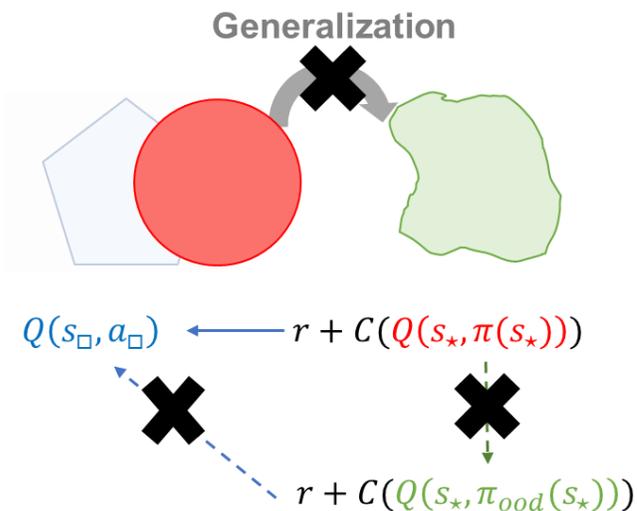
 Behavioral Policy Distribution

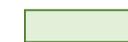
**Over-Inhibition  
when Distributions Overlap**



 Learning Policy Distribution

**Policy-OOD  
Generalization Inhibition**



 Learned OOD Policy Distribution

# Results

Table 1: Results of different algorithms and the ones equipped with RD

DATASET	TD3-N-UNC	TD3-N-UNC +RD	SAC-N-UNC	SAC-N-UNC +RD	TD3BC	TD3BC +RD	CQL	CQL +RD
HALFCHEETAH-M	66.8 ± 0.5	<b>66.8 ± 1.2</b>	65.9 ± 1.0	<b>65.9 ± 1.9</b>	48.0 ± 0.3	<b>48.3 ± 0.5</b>	47.1 ± 0.2	<b>53.0 ± 0.5</b>
HALFCHEETAH-MR	53.4 ± 3.9	<b>57.7 ± 0.9</b>	53.2 ± 5.4	<b>61.5 ± 1.4</b>	44.6 ± 0.3	44.6 ± 0.5	45.2 ± 0.6	<b>51.6 ± 0.9</b>
HALFCHEETAH-ME	97.7 ± 2.2	<b>101.1 ± 0.4</b>	99.4 ± 2.5	<b>102.5 ± 1.8</b>	90.5 ± 6.6	<b>93.9 ± 2.9</b>	81.1 ± 6.0	<b>90.2 ± 5.8</b>
HOPPER-M	41.9 ± 50.5	<b>103.0 ± 0.8</b>	45.7 ± 41.0	<b>102.8 ± 0.2</b>	60.4 ± 4.0	<b>61.0 ± 2.6</b>	65.0 ± 6.1	<b>74.9 ± 7.1</b>
HOPPER-MR	92.5 ± 18.1	<b>104.1 ± 0.8</b>	104.7 ± 0.9	104.6 ± 0.4	61.2 ± 20.5	<b>72.1 ± 8.4</b>	87.7 ± 14.4	<b>100.3 ± 3.2</b>
HOPPER-ME	100.3 ± 22.6	<b>110.7 ± 0.6</b>	110.9 ± 0.2	110.6 ± 0.3	105.4 ± 6.1	104.8 ± 2.8	93.9 ± 14.3	<b>98.2 ± 9.7</b>
WALKER2D-M	69.9 ± 35.2	<b>97.6 ± 3.4</b>	24.2 ± 28.2	<b>92.3 ± 1.3</b>	82.7 ± 5.5	<b>83.7 ± 2.7</b>	80.4 ± 3.5	<b>84.5 ± 1.0</b>
WALKER2D-MR	91.6 ± 2.7	<b>92.1 ± 2.7</b>	85.2 ± 2.7	<b>86.9 ± 3.1</b>	82.1 ± 2.5	<b>84.8 ± 1.4</b>	79.2 ± 5.0	<b>94.4 ± 2.5</b>
WALKER2D-ME	90.6 ± 45.0	<b>118.8 ± 1.2</b>	113.1 ± 9.6	<b>116.4 ± 1.5</b>	110.2 ± 0.5	110.1 ± 0.5	109.7 ± 0.5	<b>113.0 ± 0.5</b>

Table 2: Average normalized scores of our methods and previous methods on the D4RL benchmark

DATASET	BC	DT	TD3BC	CQL	IQL	EDAC	DIFFUSION-QL	SAC-N-UNC +RD	TD3-N-UNC +RD
HALFCHEETAH-R	2.2	2.2	11.0	31.3	13.7	28.4	22.0	25.4	31.0
HOPPER-R	3.7	5.4	8.4	5.3	8.4	25.3	18.3	<b>31.6</b>	<b>31.7</b>
WALKER2D-R	1.3	2.2	1.7	5.4	5.9	16.6	5.5	<b>21.2</b>	<b>21.7</b>
HALFCHEETAH-M	42.6	42.6	48.0	47.1	47.4	65.9	51.5	<b>65.9</b>	<b>66.8</b>
HOPPER-M	52.9	67.6	60.4	65.0	66.3	101.6	96.6	<b>102.8</b>	<b>103.0</b>
WALKER2D-M	75.3	74.0	82.7	80.4	78.3	92.5	87.3	92.3	<b>97.6</b>
HALFCHEETAH-MR	36.6	36.6	44.6	45.2	44.2	61.3	48.3	<b>61.5</b>	57.7
HOPPER-MR	18.1	82.7	61.2	87.7	94.7	101.0	102.0	<b>104.6</b>	<b>104.1</b>
WALKER2D-MR	26.0	66.6	82.1	79.2	73.9	87.1	98.0	86.9	92.1
HALFCHEETAH-ME	55.2	86.8	90.5	81.1	86.7	106.3	97.2	102.5	101.1
HOPPER-ME	52.5	107.6	105.4	93.9	91.5	110.7	112.3	110.6	110.7
WALKER2D-ME	107.5	108.1	110.2	109.7	109.6	114.7	111.2	<b>116.4</b>	<b>118.8</b>
HALFCHEETAH-E	91.8	87.7	96.7	97.3	94.9	106.8	96.3	<b>108.8</b>	103.1
HOPPER-E	107.7	94.2	107.8	106.5	108.8	110.1	102.6	109.8	108.8
WALKER2D-E	108.7	108.3	110.2	109.3	109.7	115.1	109.5	112.3	111.2
MUJoCo TOTAL	782.1	972.6	1020.9	1044.4	1034.0	1243.4	1158.6	<b>1252.6</b>	<b>1259.4</b>
PEN-HUMAN	25.8	73.9	-1.9	35.2	71.5	52.1	75.7	61.1	<b>77.9</b>
PEN-CLONED	38.3	67.3	9.6	27.2	37.3	68.2	60.8	53.0	65.5
ADROIT TOTAL	64.1	141.2	7.7	62.4	108.8	120.3	136.5	114.1	<b>143.4</b>
TOTAL	846.2	1113.8	1028.6	1106.8	1142.8	1363.7	1295.1	<b>1366.7</b>	<b>1402.8</b>

# Thanks

Welcome to communicate and cooperate with **Tianjin University Deep Reinforcement Learning Lab**

- Homepage: <http://rl.beiyang.ren/>