

# Implicit Q Learning: Improvements on Antmaze

Reporter:

廖修誼 (111652017)

吳泓諺 (111652040)

王裕昕 (111550066)

# Outline

- Part I (Introduction)
- Part II (Interesting findings)
- Part III (Tech. 1 - Distribution model)
- Part IV (Tech. 2 - D2RL)
- Part V (Tech. 3 - Bonus Reward)
- Part VI (Conclusion)



# Part I (Introduction)



# IQL Setting

- Offline RL
- main insight: no need to evaluate OOD actions
- method: approx. an upper expectile of distribution
- Goal: minimizing the deviation from the behavior policy

# IQL algorithm

## Algorithm 1 Implicit Q-learning

Initialize parameters  $\psi, \theta, \hat{\theta}, \phi$ .

TD learning (IQL):

**for** each gradient step **do**

$$\psi \leftarrow \psi - \lambda_V \nabla_{\psi} L_V(\psi)$$

$$\theta \leftarrow \theta - \lambda_Q \nabla_{\theta} L_Q(\theta)$$

$$\hat{\theta} \leftarrow (1 - \alpha) \hat{\theta} + \alpha \theta$$

**end for**

Policy extraction (AWR):

**for** each gradient step **do**

$$\phi \leftarrow \phi - \lambda_{\pi} \nabla_{\phi} L_{\pi}(\phi)$$

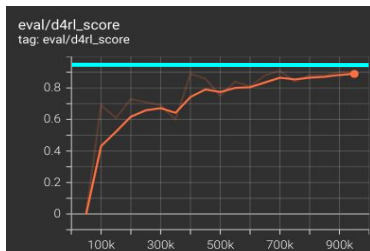
**end for**

$$L_V(\psi) = \mathbb{E}_{(s,a) \sim \mathcal{D}} [L_2^{\tau} (Q_{\hat{\theta}}(s, a) - V_{\psi}(s))]$$

$$L_Q(\theta) = \mathbb{E}_{(s,a,s') \sim \mathcal{D}} [(r(s, a) + \gamma V_{\psi}(s') - Q_{\theta}(s, a))^2]$$

$$L_{\pi}(\phi) = \mathbb{E}_{(s,a) \sim \mathcal{D}} [\exp(\beta (Q_{\hat{\theta}}(s, a) - V_{\psi}(s))) \log \pi_{\phi}(a | s)]$$

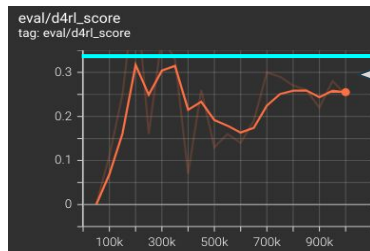
# Result Reproduce (antmaze-v0)



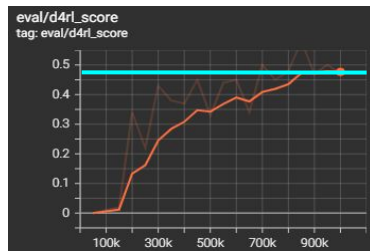
umaze,  $\tau = 0.9$



medium,  $\tau = 0.9$

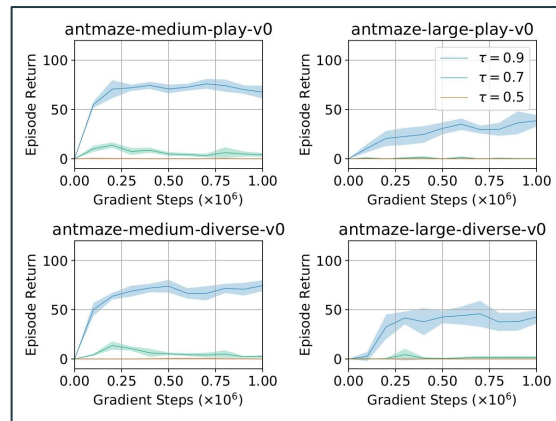


large,  $\tau = 0.9$



large-diverse,  $\tau = 0.9$

loss hyper  
parameter



# Experiment

Dataset	BC	10%BC	DT	AWAC	Onestep RL	TD3+BC	CQL	IQL (Ours)
halfcheetah-medium-v2	42.6	42.5	42.6	43.5	<b>48.4</b>	<b>48.3</b>	44.0	<b>47.4</b>
hopper-medium-v2	52.9	56.9	<b>67.6</b>	57.0	59.6	59.3	58.5	<b>66.3</b>
walker2d-medium-v2	75.3	75.0	74.0	72.4	<b>81.8</b>	83.7	72.5	78.3
halfcheetah-medium-replay-v2	36.6	40.6	36.6	40.5	38.1	<b>44.6</b>	<b>45.5</b>	<b>44.2</b>
hopper-medium-replay-v2	18.1	75.9	82.7	37.2	<b>97.5</b>	60.9	<b>95.0</b>	<b>94.7</b>
walker2d-medium-replay-v2	26.0	62.5	66.6	27.0	49.5	<b>81.8</b>	77.2	73.9
halfcheetah-medium-expert-v2	55.2	<b>92.9</b>	86.8	42.8	<b>93.4</b>	<b>90.7</b>	<b>91.6</b>	86.7
hopper-medium-expert-v2	52.5	<b>110.9</b>	<b>107.6</b>	55.8	103.3	98.0	<b>105.4</b>	91.5
walker2d-medium-expert-v2	<b>107.5</b>	<b>109.0</b>	<b>108.1</b>	74.5	<b>113.0</b>	<b>110.1</b>	<b>108.8</b>	<b>109.6</b>
locomotion-v2 total	466.7	<b>666.2</b>	<b>672.6</b>	450.7	<b>684.6</b>	<b>677.4</b>	<b>698.5</b>	<b>692.4</b>
antmaze-umaze-v0	54.6	62.8	59.2	56.7	64.3	78.6	74.0	<b>87.5</b>
antmaze-umaze-diverse-v0	45.6	50.2	53.0	49.3	60.7	71.4	<b>84.0</b>	62.2
antmaze-medium-play-v0	0.0	5.4	0.0	0.0	0.3	10.6	61.2	<b>71.2</b>
antmaze-medium-diverse-v0	0.0	9.8	0.0	0.7	0.0	3.0	53.7	<b>70.0</b>
antmaze-large-play-v0	0.0	0.0	0.0	0.0	0.0	0.2	15.8	<b>39.6</b>
antmaze-large-diverse-v0	0.0	6.0	0.0	1.0	0.0	0.0	14.9	<b>47.5</b>
antmaze-v0 total	100.2	134.2	112.2	107.7	125.3	163.8	303.6	<b>378.0</b>
total	566.9	800.4	784.8	558.4	809.9	841.2	1002.1	<b>1070.4</b>
kitchen-v0 total	<b>154.5</b>	-	-	-	-	-	144.6	<b>159.8</b>
adroit-v0 total	104.5	-	-	-	-	-	93.6	<b>118.1</b>
total+kitchen+adroit	825.9	-	-	-	-	-	1240.3	<b>1348.3</b>
runtime	10m	10m	960m	20m	$\approx 20m^*$	20m	80m	20m



## Part II (Interesting findings)

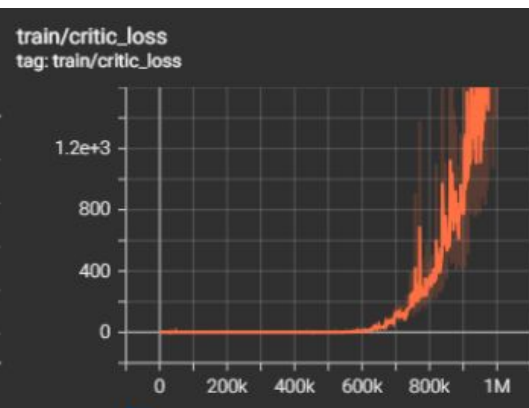
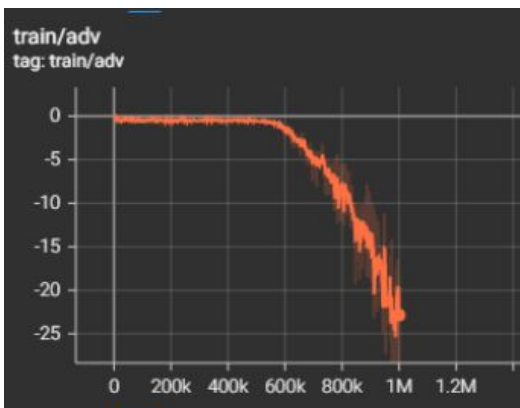
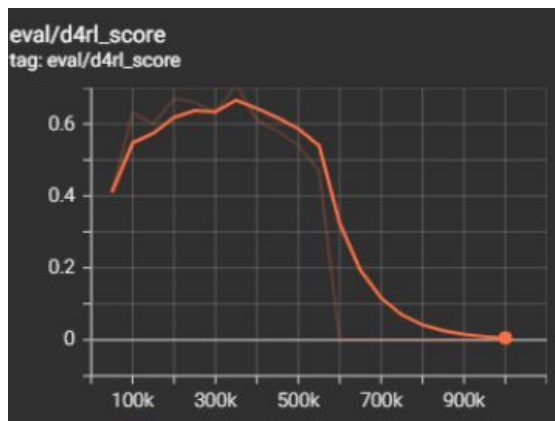




# Rethink on the equation

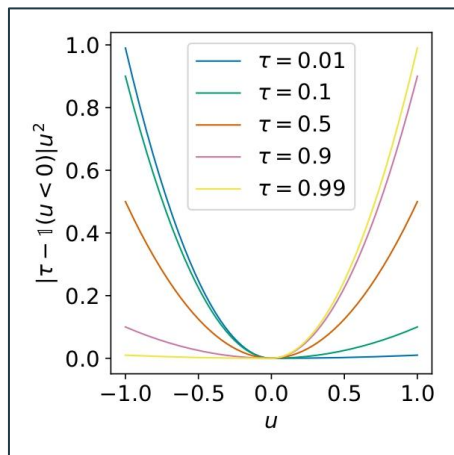
value network:	$L_V(\psi) = \mathbb{E}_{(s,a) \sim \mathcal{D}} [L_2^\tau (Q_{\hat{\theta}}(s,a) - V_\psi(s))]$	$\longrightarrow$	$Q < V$
Q network:	$L_Q(\theta) = \mathbb{E}_{(s,a,s') \sim \mathcal{D}} [(r(s,a) + \gamma V_\psi(s') - Q_\theta(s,a))^2]$	$\longrightarrow$	$Q \rightarrow r + \gamma * V$
Policy:	$L_\pi(\phi) = \mathbb{E}_{(s,a) \sim \mathcal{D}} [\exp(\beta (Q_{\hat{\theta}}(s,a) - V_\psi(s))) \log \pi_\phi(a   s)]$		

# What if there is no double Q in implementation



$$\text{adv} = q - v$$

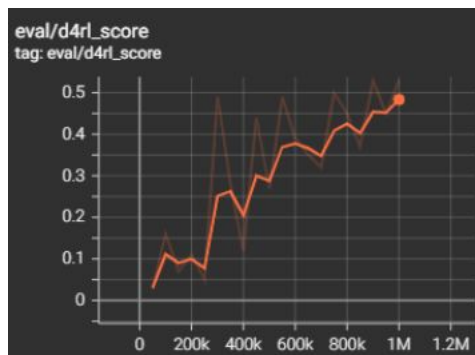
# Meaning of $\tau$



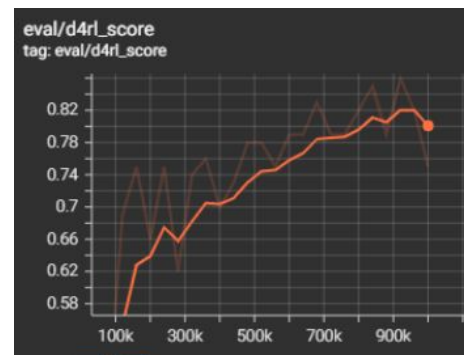
$$L_2^\tau(u) = |\tau - 1(u < 0)|u^2$$

**Lemma 2.** For all  $s$ ,  $\tau_1$  and  $\tau_2$  such that  $\tau_1 < \tau_2$  we get

$$V_{\tau_1}(s) \leq V_{\tau_2}(s).$$



expectile = 0.8



expectile = 0.9

# Evolve from quantile regression loss

## The Quantile Regression Loss

- Given that the derivative of  $L(x; Z, \tau)$  is  $F_Z(x) - \tau$ , we can recover the QR loss by integration

Quantile regression (QR) loss:

$$L_{QR}(x; Z, \tau) = (\tau - 1) \int_{-\infty}^x (z - x) dF_Z(z) + \tau \int_x^{\infty} (z - x) dF_Z(z)$$

(It is easy to verify that  $\frac{d}{dx} L_{QR}(x; Z, \tau) = F_Z(x) - \tau$  by the Leibniz integral rule)

Alternative expression of QR loss:

$$\begin{aligned} \frac{d}{dx} \left( \int_{a(x)}^{b(x)} f(x, t) dt \right) \\ = f(x, b(x)) \cdot \frac{d}{dx} b(x) - f(x, a(x)) \cdot \frac{d}{dx} a(x) + \int_{a(x)}^{b(x)} \frac{\partial}{\partial x} f(x, t) dt \end{aligned}$$

$$\rho_{\tau}(y) := y(\tau - \mathbb{1}\{y < 0\})$$

$$L_{QR}(x; Z, \tau) = E_Z[\rho_{\tau}(Z - x)]$$

## Value loss

$$L_V(\psi) = \mathbb{E}_{(s,a) \sim \mathcal{D}} [L_2^{\tau}(Q_{\hat{\theta}}(s, a) - V_{\psi}(s))]$$

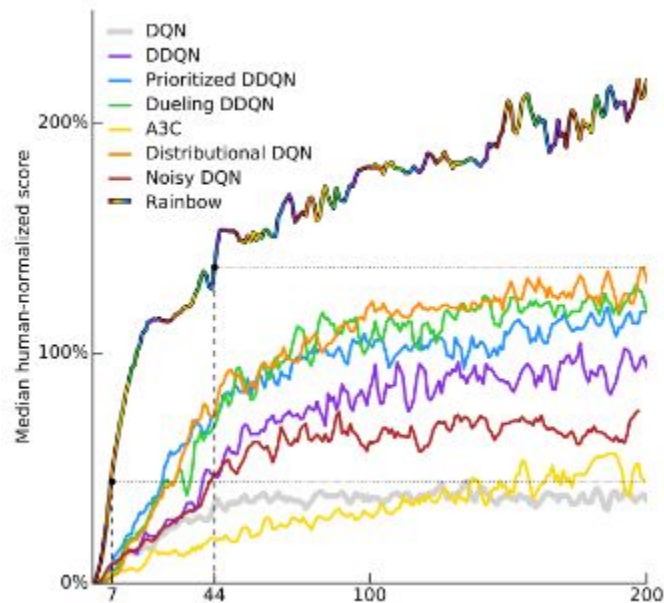
$$L_2^{\tau}(u) = |\tau - \mathbb{1}(u < 0)|u^2.$$



## Part III (Tech. 1 - Distribution model)



# Motivation idea



improve performance with  
Distributionalizing IQL

# Comparison : our method & quantile

## implicit quantile network

$$\delta_t^{\tau, \tau'} = r_t + \gamma Z_{\tau'}(x_{t+1}, \pi_{\beta}(x_{t+1})) - Z_{\tau}(x_t, a_t). \quad (2)$$

Then, the IQN loss function is given by

$$\mathcal{L}(x_t, a_t, r_t, x_{t+1}) = \frac{1}{N'} \sum_{i=1}^N \sum_{j=1}^{N'} \rho_{\tau_i}^{\kappa} \left( \delta_t^{\tau_i, \tau'_j} \right), \quad (3)$$

## Distributional IQL

$$L_V(\psi) = \mathbb{E}_{(s,a) \sim \mathcal{D}} \left[ L_2^{\tau} (Q_{\hat{\theta}}(s, a) - V_{\psi}(s)) \right]$$

$$L_Q(\theta) = \mathbb{E}_{(s,a,s') \sim \mathcal{D}} \left[ (r(s, a) + \gamma V_{\psi}(s') - Q_{\theta}(s, a))^2 \right]$$

$$\text{mean\_loss} = \frac{1}{N} \sum_{i=1}^N \sum_{j=1}^T \sum_{k=1}^T (q_{ij} - v_{ik})^2 \cdot |\tau - 1_{\{q_{ij} - v_{ik} < 0\}}| \cdot \text{prob}_{ij} \cdot \text{prob}_{ik}$$

implementation follows from IQN

## Guess , benefit of our method : Two step policy improvement

$$L_Q(\theta) = \mathbb{E}_{(s,a,s') \sim \mathcal{D}} \left[ (r(s,a) + \gamma V_\psi(s') - Q_\theta(s,a))^2 \right]$$

$$L_V(\psi) = \mathbb{E}_{(s,a) \sim \mathcal{D}} \left[ L_2^\tau (Q_{\hat{\theta}}(s,a) - V_\psi(s)) \right]$$

$$\delta = r + \gamma G_{\theta'}(\tau'; s', a') - G_\theta(\tau; s, a).$$

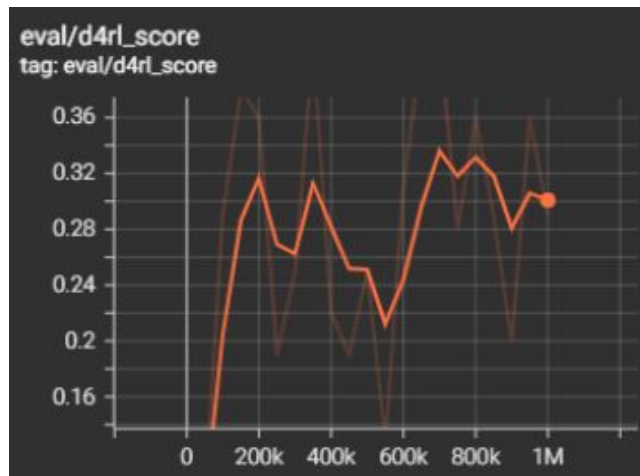
$$\mathcal{L}_\kappa(\delta; \tau) = \begin{cases} |\tau - \mathbb{1}(\delta < 0)| \cdot \delta^2 / (2\kappa) & \text{if } |\delta| \leq \kappa \\ |\tau - \mathbb{1}(\delta < 0)| \cdot (|\delta| - \kappa / 2) & \text{otherwise .} \end{cases}$$



# However, experiment result is ...



**antmaze-medium-play-v0 – expectile 0.9**



**antmaze-large-play-v0, expectile 0.9**



## Part IV (Tech. 2 - D2RL)



# Motivation idea

- The problem of choosing architecture designs has been largely ignored.
- Information loss when forwarding through layers.
- The effective rank of the feature matrix is low.
- Add skip connections from the input.

# Skip Connection

original:

```
nn.Linear(in_dim, hidden_dim)
```

```
nn.Linear(hidden_dim, hidden_dim)
```

```
nn.Linear(hidden_dim, out_dim)
```

with skip connections:

```
nn.Linear(in_dim, hidden_dim)
```

```
nn.Linear(in_dim + hidden_dim, hidden_dim)
```

```
nn.Linear(hidden_dim, out_dim)
```

## Results — srnk

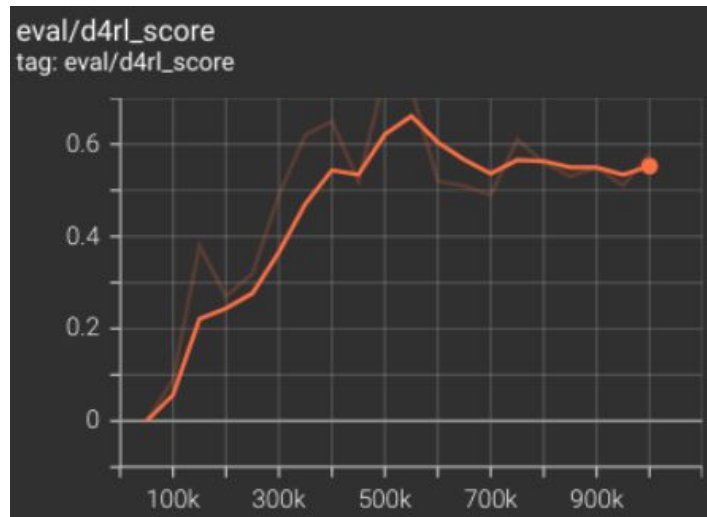
$$\text{srnk}_\delta(\Phi) = \min\{k : \frac{\sum_{i=1}^k \sigma_i(\Phi)}{\sum_{i=1}^d \sigma_i(\Phi)} \geq 1 - \delta\}$$

	antmaze-large		halfcheetah-expert	
1M-steps	IQL	IQL+SC	IQL	IQL+SC
Policy	227	<b>232</b>	226	<b>232</b>
Q-network	223	<b>231</b>	225	<b>233</b>

# Results

- competitive performance similar to original IQL
- outperforms IQL on antmaze-large
- prevent the reduction of effective ranks

better convergence  
on antmaze-large





## Part V (Tech. 3 - Bonus Reward)



# Motivation idea

## Offline Reinforcement Learning as Anti-Exploration

Shideh Rezaeifar<sup>\*1</sup>, Robert Dadashi<sup>\*2</sup>, Nino Vieillard<sup>2,3</sup>, Léonard Hussenot<sup>2,4</sup>, Olivier Bachem<sup>2</sup>, Olivier Pietquin<sup>2</sup>, and Matthieu Geist<sup>2</sup>

<sup>1</sup>University of Geneva

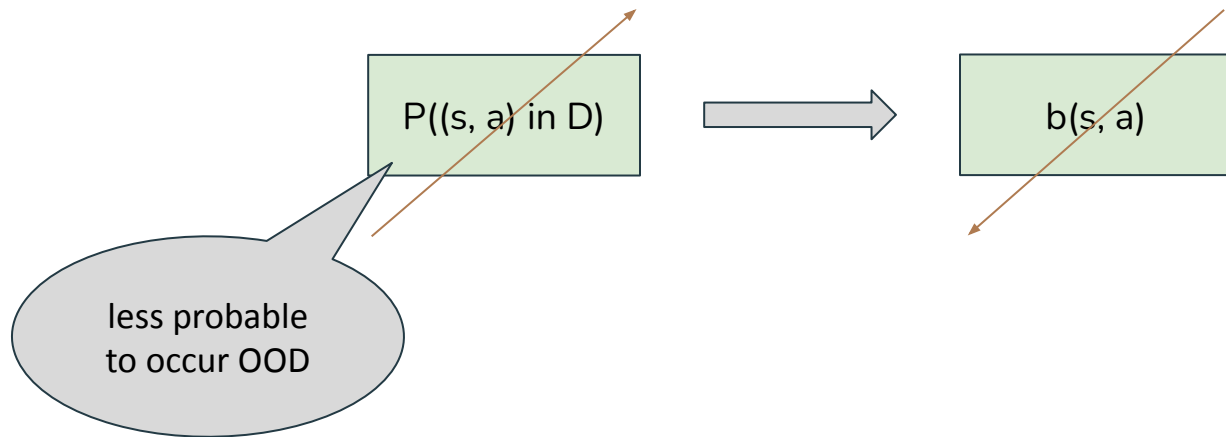
<sup>2</sup>Google Research, Brain Team

<sup>3</sup>Univ. Lorraine, CNRS, Inria, IECL, F-54000 Nancy, France

<sup>4</sup>Univ. Lille, CNRS, Inria Scool, UMR 9189 CRISAL



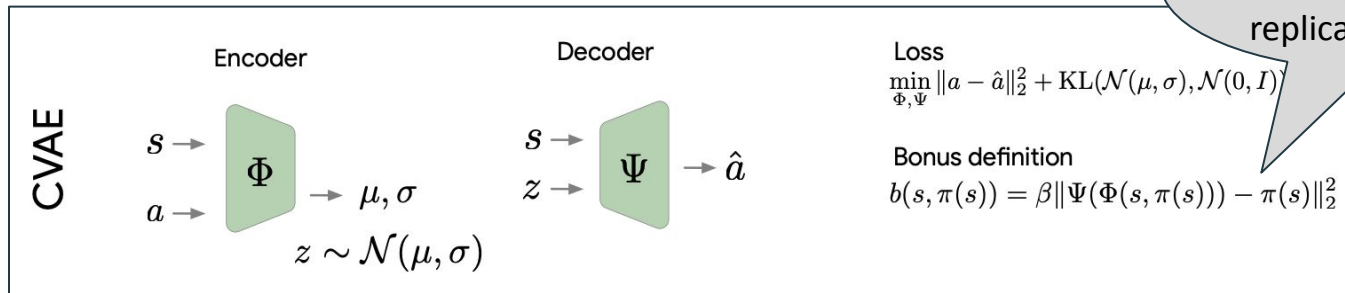
# (subtractive) Bonus Reward Property



# Train Bonus Reward (CVAE)

## Algorithm 1 CVAE training.

- 1: Initialize CVAE networks  $\Phi$  and  $\Psi$
- 2: **for** step  $i = 0$  to  $N$  **do**
- 3:   Sample a minibatch of  $k$  state-action pairs  $\{(s_t, a_t), t = 1, \dots, k\}$  from  $\mathcal{D}$
- 4:   Train  $\Phi$  and  $\Psi$  using  $\mathcal{L}_{\Phi, \Psi}$ , see Eq. (5)

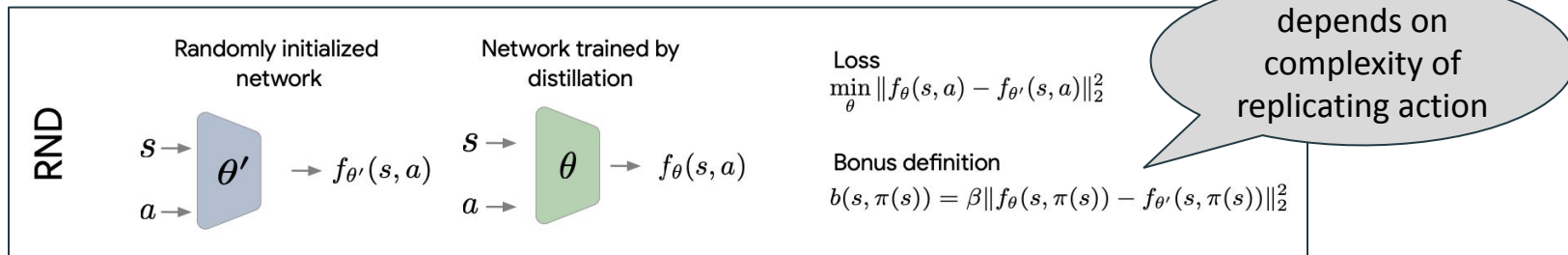


depends on  
complexity of  
replicating action

# Train Bonus Reward (RND)

## Algorithm 1 CVAE training.

- 1: Initialize CVAE networks  $\Phi$  and  $\Psi$
- 2: **for** step  $i = 0$  to  $N$  **do**
- 3:   Sample a minibatch of  $k$  state-action pairs  $\{(s_t, a_t), t = 1, \dots, k\}$  from  $\mathcal{D}$
- 4:   Train  $\Phi$  and  $\Psi$  using  $\mathcal{L}_{\Phi, \Psi}$ , see Eq. (5)



# Apply Bonus Reward to IQL

## Algorithm 1 Implicit Q-learning

Initialize parameters  $\psi, \theta, \hat{\theta}, \phi$ .

TD learning (IQL):

**for** each gradient step **do**

$$\psi \leftarrow \psi - \lambda_V \nabla_{\psi} L_V(\psi)$$

$$\theta \leftarrow \theta - \lambda_Q \nabla_{\theta} L_Q(\theta)$$

$$\hat{\theta} \leftarrow (1 - \alpha) \hat{\theta} + \alpha \theta$$

**end for**

Policy extraction (AWR):

**for** each gradient step **do**

$$\phi \leftarrow \phi - \lambda_{\pi} \nabla_{\phi} L_{\pi}(\phi)$$

**end for**

$$L_V(\psi) = \mathbb{E}_{(s,a) \sim \mathcal{D}} [L_2^T (\cancel{Q_{\hat{\theta}}(s,a)} - V_{\psi}(s))]$$

$$L_V(\psi) = \mathbb{E}_{(s,a) \sim \mathcal{D}} [L_2^T ((Q_{\hat{\theta}}(s,a) - b(s,a)) - V_{\psi}(s))]$$

$$L_Q(\theta) = \mathbb{E}_{(s,a,s') \sim \mathcal{D}} [(r(s,a) + \gamma \cancel{V_{\psi}(s')} - Q_{\theta}(s,a))^2]$$

$$L_Q(\theta) = \mathbb{E}_{(s,a,s') \sim \mathcal{D}} [\mathbb{E}_{a' = \mu_{\theta}(s'), \epsilon \sim N(0, \sigma^2 I)} [(r(s,a) + \gamma (Q_{\psi}(s', a' + \epsilon) - b(s', a')) - Q_{\theta}(s,a))^2]$$

$$L_{\pi}(\phi) = \mathbb{E}_{(s,a) \sim \mathcal{D}} [\exp(\beta (\cancel{Q_{\hat{\theta}}(s,a)} - V_{\psi}(s))) \log \pi_{\phi}(a | s)]$$

$$L_{\pi}(\phi) = \mathbb{E}_{(s,a) \sim \mathcal{D}} [\exp(\beta (Q_{\hat{\theta}}(s,a) - b(s,a)) - V_{\psi}(s)) \log \pi_{\phi}(a | s)]$$



# Part VI (Conclusion)



# Recall

- Part I (Introduction)
- Part II (interesting findings)
- Part III (Tech. 1 - Distribution model): modify update formulas to distribution form
- Part IV (Tech. 2 - D2RL): modify neural layers
- Part V (Tech. 3 - bonus reward): modify reward

**Thank you for your attention**