# Optimistic Critic Reconstruction and Constrained Fine-Tuning for General Offline-to-Online RL
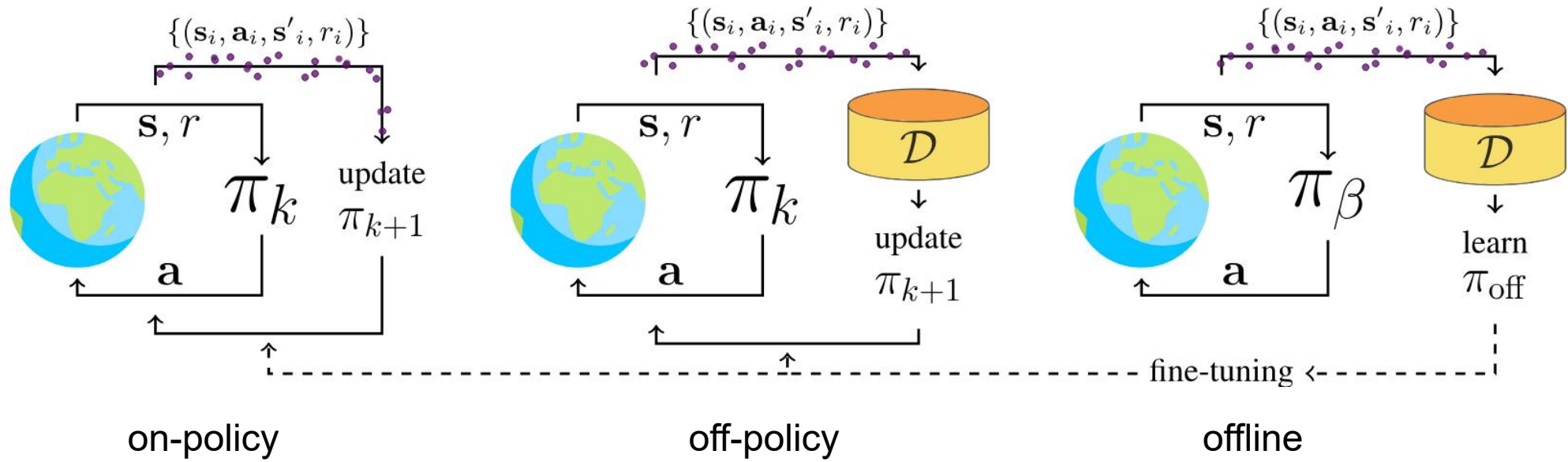
Qin-Wen Luo*[1] , Ming-Kun Xie*[1,2], Ye-Wen Wang*[1] , Sheng-Jun Huang✉[1]

[1] Nanjing University of Aeronautics and Astronautics, Nanjing, China

[2] RIKEN Center for Advanced Intelligence Project, Tokyo, Japan

# Offline-to-online RL



$\{(\mathbf{s}_i, \mathbf{a}_i, \mathbf{s}'_i, r_i)\}$

$\mathbf{s}, r$

$\pi_k$ update $\pi_{k+1}$

$\mathbf{a}$

$\{(\mathbf{s}_i, \mathbf{a}_i, \mathbf{s}'_i, r_i)\}$

$\mathbf{s}, r$

$\pi_k$ $\mathcal{D}$

$\mathbf{a}$ update $\pi_{k+1}$

$\{(\mathbf{s}_i, \mathbf{a}_i, \mathbf{s}'_i, r_i)\}$

$\mathbf{s}, r$

$\pi_\beta$ $\mathcal{D}$

$\mathbf{a}$ learn $\pi_{\text{off}}$

fine-tuning

on-policy                    off-policy                    offline

combat distribution shift          policy constraint/value regularization

maintain pessimism

hinder efficiency/individually designed

# Two mismatches

$$J(\pi) = \mathbb{E}_\pi \left[ \sum_{t=0}^{\infty} \gamma^t r(s_t, a_t) \right]$$

$$Q^{k+1}(s,a) = r(s,a) + \mathbb{E}_{s' \sim P(\cdot|s,a), a' \sim \pi_k(\cdot|s')} \left[ Q^k(s', a') \right] \quad \text{policy evaluation}$$

$$\pi_{k+1} = \arg\max_\pi Q^{\pi_k}(s,a) \quad \text{policy improvement}$$

evaluation mismatch

$$J(\pi) = \mathbb{E}_\pi \left[ \sum_{t=0}^{\infty} \gamma^t \left( r(s_t, a_t) - \alpha \cdot f\left( \frac{\pi(a_t|s_t)}{\mu(a_t|s_t)} \right) \right) \right]$$
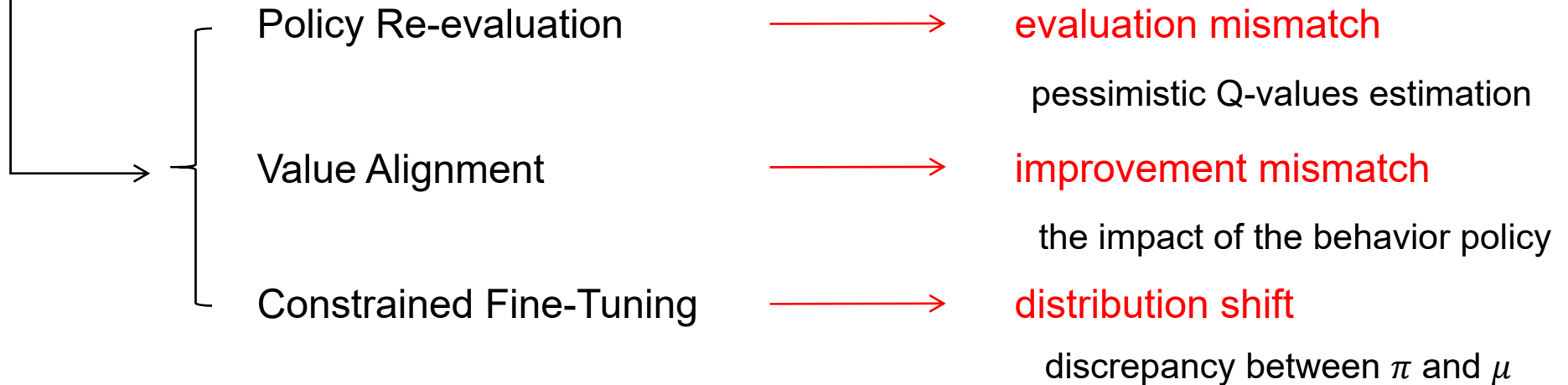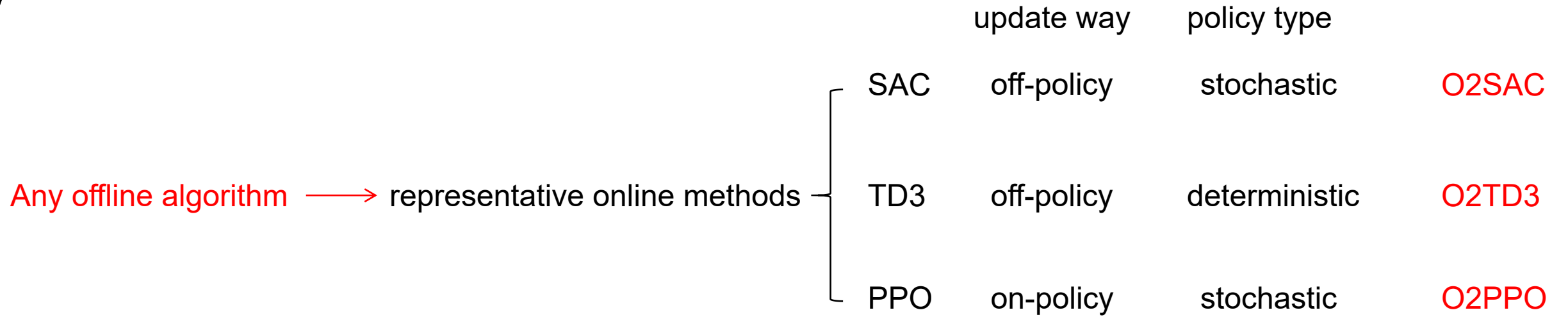
improvement mismatch

$$Q^{k+1}(s,a) = r(s,a) + \mathbb{E}_{s' \sim P(\cdot|s,a), a' \sim \pi_k(\cdot|s')} \left[ Q^k(s', a') - \alpha \cdot f\left( \frac{\pi_k(a'|s')}{\mu(a'|s')} \right) \right]$$

$$\pi_{k+1} = \arg\max_\pi Q^{\pi_k}(s,a) - \alpha \cdot f\left( \frac{\pi_k(a|s)}{\mu(a|s)} \right)$$

# Method

# Policy Re-evaluation

off-policy evaluation (OPE)

Fitted Q-learning (FQE) $\longrightarrow$ O2SAC, O2TD3    consistent with online update

theoretical guarantee of FQE

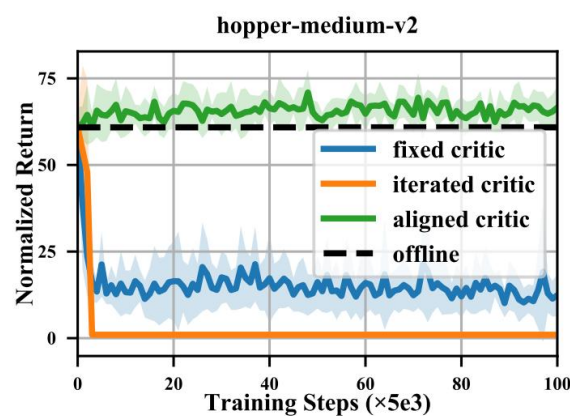$$\max_{(s,a)\in S\times A} \frac{d^{\pi_\theta}(s,a)}{d^{\mu}(s,a)} \leq C.$$

$$|Q^\pi - \hat{Q}^\pi| \leq \frac{1-\gamma^K}{1-\gamma}\sqrt{C\epsilon} + \gamma^K\bar{V}, \; where \; \epsilon := \frac{22\bar{V}^2\log(|\mathcal{F}|/\delta)}{|\mathcal{D}|} + 20d^\pi_F.$$

optimistic critic

Fitted returns $\longrightarrow$ O2PPO    the fitted retrurns only approximate $V^\mu(s)$

# Value Alignment
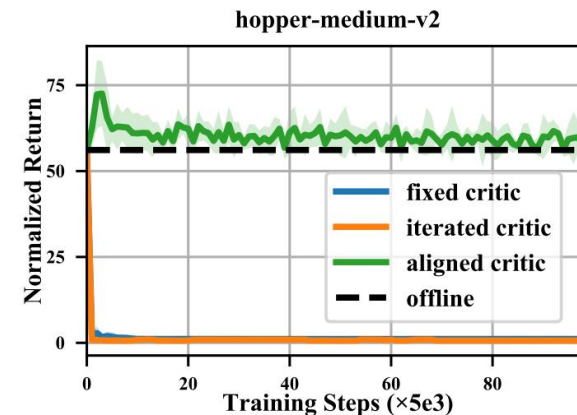
inconsistent offline update ⎫
                           ⎬ improvement mismatch
imperfect data for FQE     ⎭



(a) The performance of SAC    (b) The performance of TD3

O2SAC:  $Q(s,a) = V(s) + \alpha \log \pi(a|s)$ ⟹ $Q'_\mu(s,a) = \min\left(Q_{\bar{\mu}}(s,\dot{a}) - \alpha(\log \pi_{\text{off}}(\dot{a}|s) - \log \pi_{\text{off}}(a|s)), Q_{\bar{\mu}}(s,a)\right)$

value bound   $V_{fqe}(s) \leq V_{align}(s) \leq V_{\dot{a}}(s)$

O2TD3:  $Q(s,a)/Q(s,\dot{a}) \sim N(\dot{a}, \Sigma)$ ⟹ $Q'(s,a) = \min\left(Q_{\bar{\mu}}(s,\tilde{a}), \dfrac{Q(s,\dot{a})}{1 + k \cdot \max(d(a,\dot{a})^2, \sigma^2)}\right)$

O2PPO:  $A_\alpha(s,a) = \alpha \log \pi_{\text{off}}(a|s) + \alpha \mathcal{H}(\pi_{\text{off}}(\cdot|s))$ ⟹ $A'(s,a) = A(s,a) + \beta A_\alpha(s,a)$

equivalent constraint   $-\mathbb{E}_{s \sim R} \mathbb{E}_{a \sim \pi_{\theta_k}(\cdot|s)} \left[\dfrac{\pi_\theta}{\pi_{\theta_k}} A_\alpha(s,a)\right] \iff CELoss(\pi_\theta, \pi_{\text{ref}}) + C$

# Constrained Fine-Tuning

Constrained MDP

$$\max \mathbb{E}_{\pi}\left[\sum_{t=0}^{\infty} \gamma_t r_t(s_t, a_t)\right] \quad \text{s.t. } \mathbb{E}_{\pi}\left[f(\pi(a_t|s_t), \pi_{\text{ref}}(a_t|s_t))\right] < \tau$$

$$
\begin{aligned}
L(\theta) &= \max \mathbb{E}_{\pi_\theta}\left[Q_\mu^{\pi_\theta}(s, a) - \lambda f(\pi_\theta(a|s), \pi_{\text{ref}}(a|s))\right] \\
L(\mu) &= \min \mathbb{E}_{(s,a,r,s')\sim R}\left[(Q_\mu^{\pi_\theta}(s, a) - y)^2\right] \\
y &= r + \gamma \mathbb{E}_{a'\sim\pi_\theta(\cdot|s')}\left[Q_{\bar{\mu}}^{\pi_\theta}(s', a') - \lambda f(\pi_\theta(a'|s'), \pi_{\text{ref}}(a'|s'))\right] \\
L(\lambda) &= \min_{\lambda \geq 0} -\lambda\left[\mathbb{E}_{\pi_\theta}(f(\pi_\theta(a|s), \pi_{\text{ref}}(a|s))) - \tau\right]
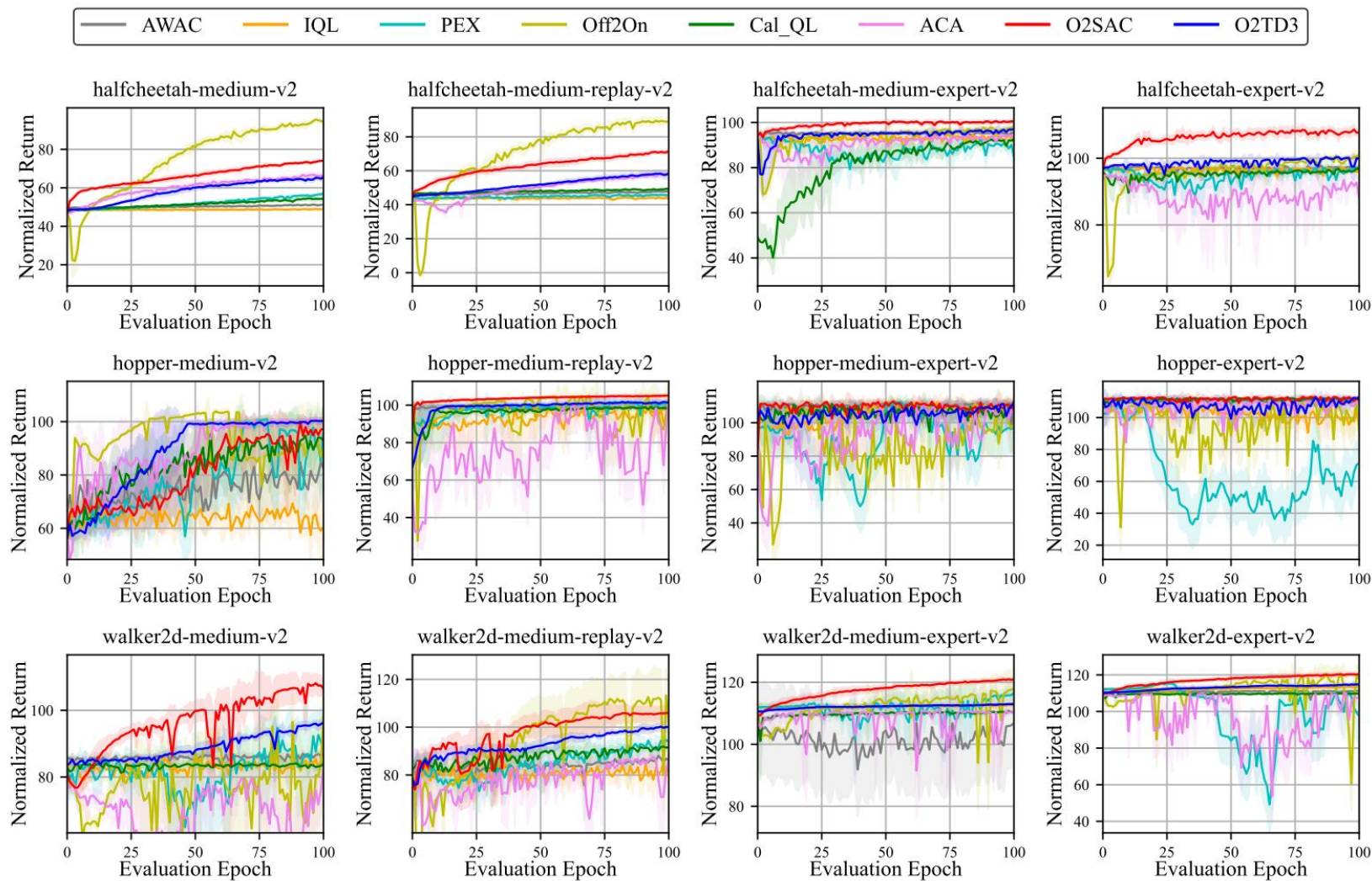\end{aligned}
\tag{20}
$$

**Corollary 4.5.** *With the penalty $f(\pi, \pi_{\text{ref}})$ defined before and appropriate learning rates, algorithm of Eq. (20) almost surely to a fixed point $(\theta^*, \mu^*, \lambda^*)$, where $\lambda^* = 0$, $\theta^*$ and $\mu^*$ are corresponded to $\pi^*$ and $Q^*$, which are optimal in the MDP without constraint.*

O2SAC: $\quad f(\mu, \phi) = D_{KL}(\mu|\phi)$

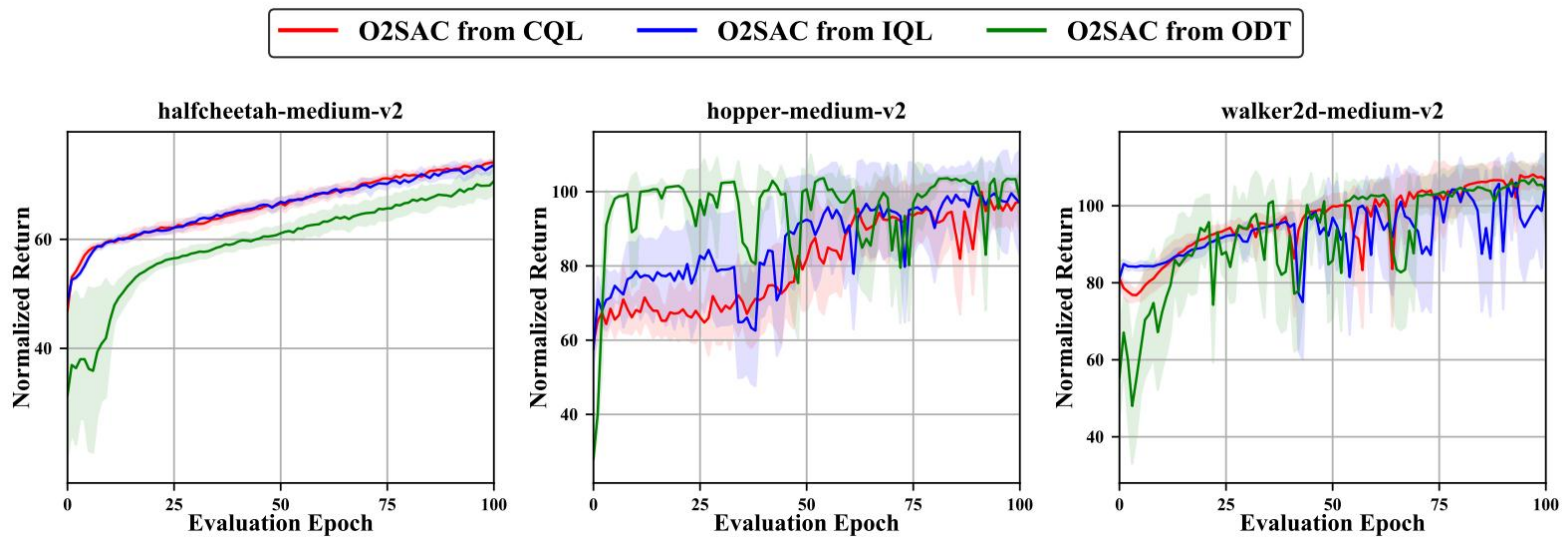O2TD3: $\quad f(\mu, \phi) = MSE(\mu, \phi)$

O2PPO: $\quad A_\alpha(s, a) = \alpha \log \pi_{\text{ref}}(a|s) + \alpha \mathcal{H}(\pi_{\text{ref}}(\cdot|s))$
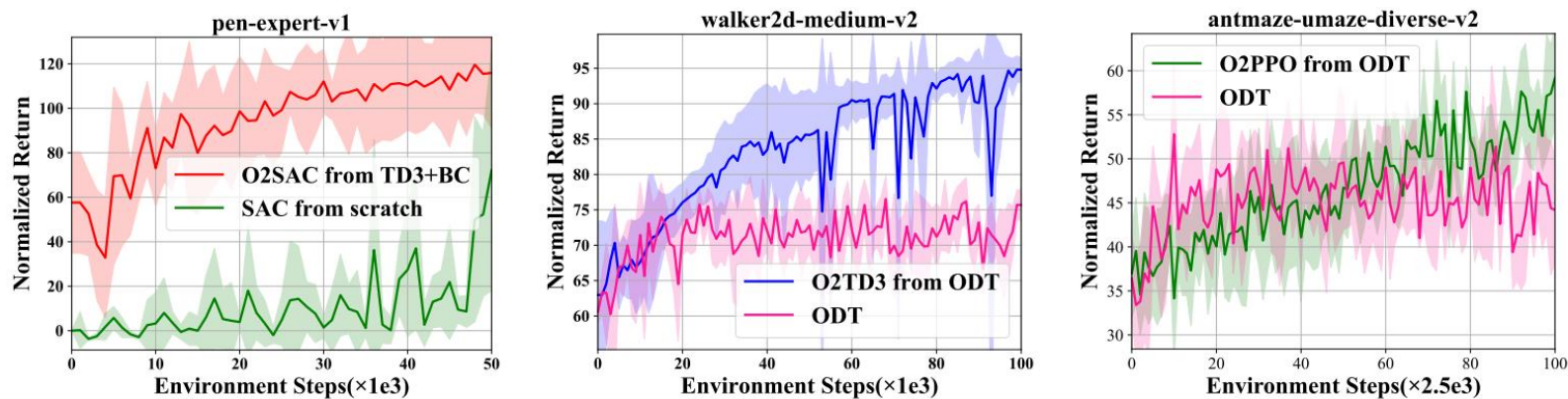
# Experiments



| Dataset | IQL | PEX | Cal-QL | O2TD3 | O2SAC | O2PPO |
|---------|-----|-----|--------|-------|-------|-------|
| U-v2 | 80.8→80.8 | 85.0→96.2 | 80.8→97.0 | 92.8→95.8 | 92.8→93.6 | 77.3→98.0 |
| U-D-v2 | 56.6→35.8 | 12.6→16.0 | 23.8→71.2 | 38.4→52.2 | 43.8→79.8 | 56.4→86.3 |
| total | 137.4→116.6 | 97.6→112.2 | 104.6→168.2 | 128.6→142.0 | 131.2→148.0 | 133.7→184.3 |

# Experiments



The fine-tuning performance achieved by initializing from different offline algorithms



(a) From TD3+BC to SAC    (b) From online DT to TD3    (c) From online DT to PPO

The fine-tuning performance achieved by transferring to three online algorithms from their heterogeneous offline algorithms.

# Code