# Score Regularized Policy Optimization through Diffusion Behavior

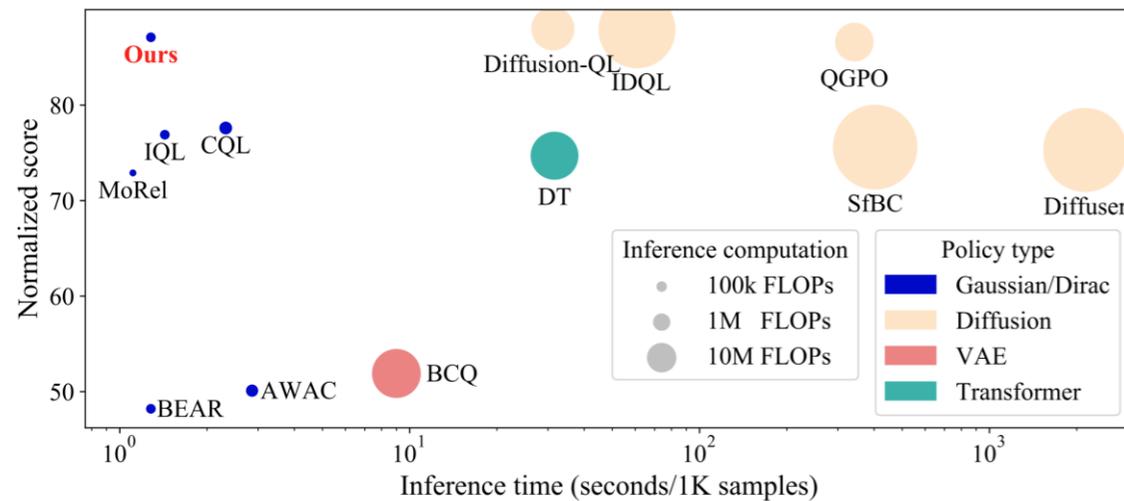Huayu Chen, Cheng Lu, Zhengyi Wang, Hang Su, Jun Zhu    **Tsinghua University**

## Overview

- We propose an offline RL algorithm for continuous control tasks.
- We leverage a powerful diffusion behavior model to regularize policy training.
- We **entirely** avoid iterative sampling from diffusion models in both during training and evaluation. This greatly increases computational efficiency



Key idea: diverse diffusion behavior but simple Dirac inference policy

## Background

$$\min_\theta \mathbb{E}_{s\sim\mathcal{D}^\mu} \underbrace{D_{\mathrm{KL}}\left[\pi^*(\cdot|s)||\pi_\theta(\cdot|s)\right]}_{\textit{Forward KL}} \Leftrightarrow \max_\theta \mathbb{E}_{(s,a)\sim\mathcal{D}^\mu}\underbrace{\left[\frac{1}{Z(s)}\log\pi_\theta(a|s)\,e^{\beta Q_\phi(s,a)}\right]}_{\textit{Weighted Regression}}, \quad (2)$$

$$\updownarrow$$

$$\min_\theta \mathbb{E}_{s\sim\mathcal{D}^\mu} \underbrace{D_{\mathrm{KL}}\left[\pi_\theta(\cdot|s)||\pi^*(\cdot|s)\right]}_{\textit{Reverse KL}} \Leftrightarrow \boxed{\max_\theta \mathbb{E}_{s\sim\mathcal{D}^\mu,a\sim\pi_\theta}Q_\phi(s,a) - \frac{1}{\beta}D_{\mathrm{KL}}\left[\pi_\theta(\cdot|s)||\mu(\cdot|s)\right]}_{\textit{Behavior-Regularized Policy Optimization}}. \quad (3)$$
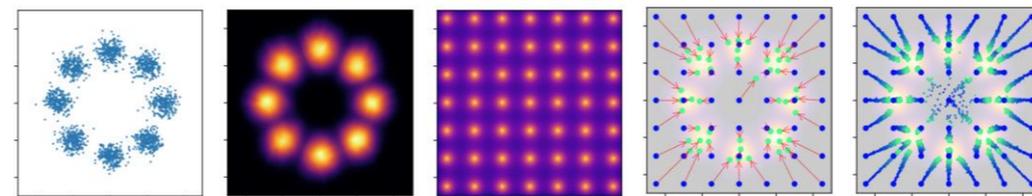
**Our initial loss function**

## Method Derivation

**Decomposing the KL term:**
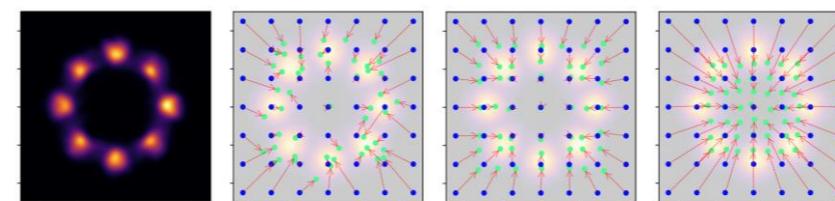
$$\mathcal{L}_\pi(\theta) = \underbrace{\mathbb{E}_{s\sim\mathcal{D}^\mu,a\sim\pi_\theta}Q_\phi(s,a)}_{\text{Policy optimization}} + \frac{1}{\beta}\underbrace{\mathbb{E}_{s\sim\mathcal{D}^\mu,a\sim\pi_\theta}\log\mu(a|s)}_{\text{Behavior regularization}} + \frac{1}{\beta}\underbrace{\mathbb{E}_{s\sim\mathcal{D}^\mu}\mathcal{H}(\pi_\theta(\cdot|s))}_{\text{Entropy (often constant[1])}}.$$

**Applying the chain rule and the reparameterization trick:**

$$\nabla_\theta\mathcal{L}_\pi(\theta) = \mathbb{E}_{s\sim\mathcal{D}^\mu}\left[\nabla_a Q_\phi(s,a)|_{a=\pi_\theta(s)} + \frac{1}{\beta}\underbrace{\nabla_a\log\mu(a|s)|_{a=\pi_\theta(s)}}\right]\nabla_\theta\pi_\theta(s).$$

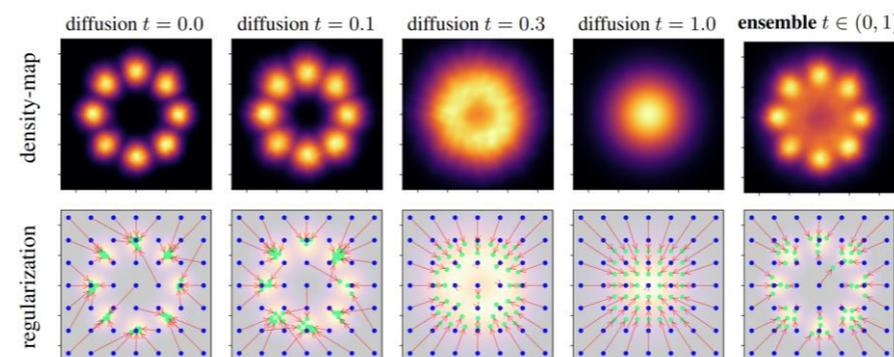can be estimated by **diffusion** models



(a) Data samples (from behavior dataset)  (b) Behavior density $\hat{\mu}(a)$ (by diffusion models)  (c) Quadratic Q-functions (stacked)  (d) Behavior regularization ($\frac{1}{\beta}=0 \rightarrow \frac{1}{\beta}=1$)  (e) Result policy shift (by varying $0 \leq \frac{1}{\beta} \leq 1$)

● w/o behavior regularization    ● w/ behavior regularization



Behavior Density (VAEs)    BCQ (Fujimoto et al., 2019)    BEAR (Kumar et al., 2019)    TD3+BC (Fujimoto & Gu, 2021)

**Ensembling various diffusion times t for better regularization result:**



diffusion t = 0.0   diffusion t = 0.1   diffusion t = 0.3   diffusion t = 1.0   **ensemble** $t \in (0,1)$

## More Experimental Reuslt