

## Appendix

### A MORE TOY EXPERIMENTS

Here we describe an additional toy experiment on a bandit task. Actions are again in a real-valued 2D space,  $\mathbf{a} \in [-1, 1]^2$ . The offline data  $\mathcal{D} = \{(\mathbf{a}_i)\}_{i=1}^{10000}$  are collected by sampling actions equally from four Gaussian distributions with centers  $\boldsymbol{\mu} \in \{(-0.8, 0.8), (0.8, 0.8), (0.8, -0.8), (-0.8, -0.8)\}$  and standard deviations  $\boldsymbol{\sigma}_d = (0.05, 0.05)$ , as depicted in the first panel of Figure 4. We conduct the same experiments as the ones in our main paper (see figure 1) and show the performance in Figure 4. The only difference in this experiment is that the samples are now in the corners of the action space.

For behavior-cloning experiments, we observe that only our diffusion model could recover the original data distribution while the prior regularization methods fail in some way. For example, CVAE could only capture the two diagonal modes and place density between them, while MMD tends to align density around the boundaries because of its Tanh-Gaussian policy. For Q-learning experiments, we observe that the prior regularization methods typically push the policy to converge to sub-optimal solutions in BCQ and BEAR while preventing the policy of TD3+BC from being concentrated on the right corner. However, the policy of Diffusion-QL successfully converges to the optimal bottom corner. The ablation study experiments are consistent with our conclusion in the main paper.

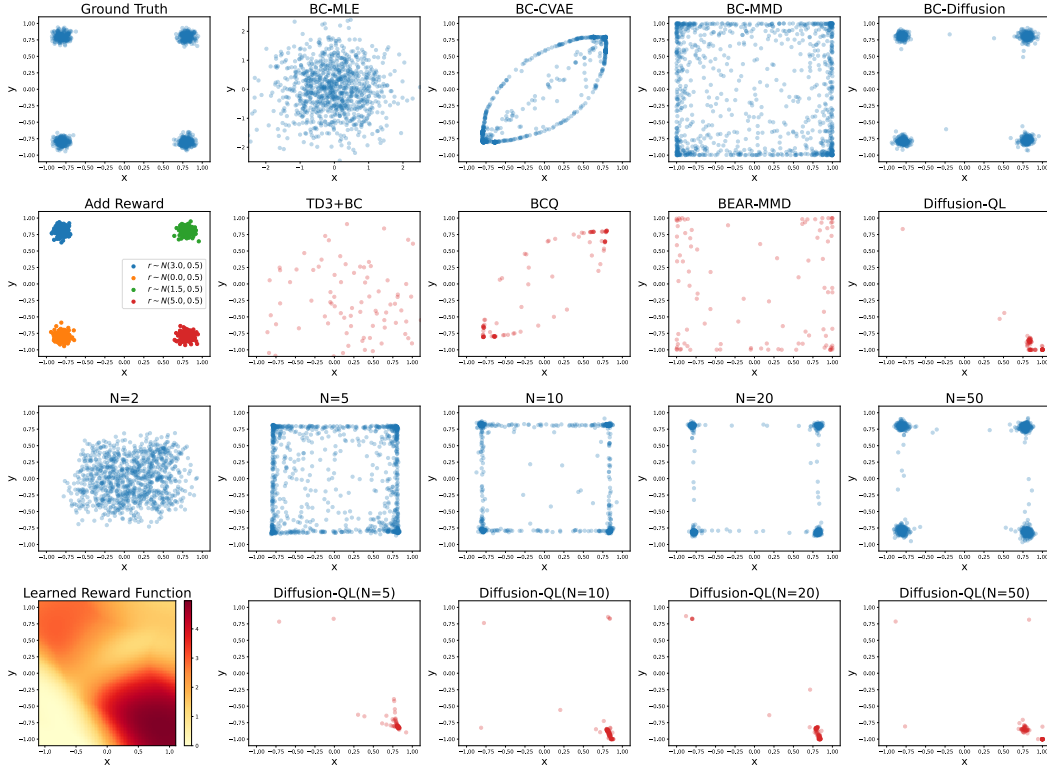


Figure 4: Bandit experiment with a strong multi-modal behavior policy. The first row shows the comparison of behavior-cloning results between our method and prior methods. The second row shows the comparison results with Q-learning involved. The third row shows the ablation study of  $N$  for BC-Diffusion. The fourth row shows the learned reward function and the ablation study of  $N$  for Diffusion-QL.

## B IMPLEMENTATION DETAILS

**Diffusion policy.** We build our policy as an MLP-based conditional diffusion model. Following the parameterization of [Ho et al. \(2020\)](#), the model itself is a residual model,  $\epsilon_\theta(\mathbf{a}^i, \mathbf{s}, i)$ , where the  $i$  is the last timestep and  $\mathbf{s}$  is the state condition. We model  $\epsilon_\theta$  as a 3-layer MLPs with Mish activations and we use 256 hidden units for all networks. The input of  $\epsilon_\theta$  is the concatenation of the last step action vector, the current state vector, and the sinusoidal positional embedding of timestep  $i$ . The output of  $\epsilon_\theta$  is the predicted residual at diffusion timestep  $i$ .

**Q networks.** We build two Q networks with the same MLP setting as our diffusion policy, which has 3-layer MLPs with Mish activations and 256 hidden units for all networks.

## C EXPERIMENTAL DETAILS

We train our algorithm with 2000 epochs for Gym tasks and 1000 epochs for the other tasks, where each epoch consists of 1000 gradient steps. For the Gym locomotion tasks, we average mean returns over 6 independently trained models and 10 trajectories per mode. For the other tasks, we average over 6 independently trained models and 100 evaluation trajectories. Following the convention of supervised learning, the best performed model is saved during training and used for final evaluation.

We use the original task rewards from MuJoCo Gym and Kitchen tasks. We standardize Adroit task rewards for training stability. We modify the rewards according to the suggestion of CQL [\(Kumar et al., 2020\)](#) for the AntMaze datasets.

## D HYPERPARAMETERS

For Diffusion-QL, we consider three hyperparameters in total: learning rate, Q-learning weight  $\eta$ , and whether to use max Q backup from CQL [\(Kumar et al., 2020\)](#). For learning rate, we consider values in the grid  $\{1 \times 10^{-3}, 3 \times 10^{-4}, 3 \times 10^{-5}\}$  for the policy, while we use a fixed learning rate,  $3 \times 10^{-4}$ , for Q-networks. For  $\eta$ , we consider values according to the characteristics of different domains, as we mentioned in the description of datasets that Adroit and Kitchen tasks require more policy regularization and AntMaze tasks require more Q-learning. For max Q backup, we only apply it on AntMaze tasks. Based on these considerations, we provide our hyperparameter setting in Table [3](#).

## E OPTIMAL RESULTS

If a small amount of online experience is provided during the evaluation stage for model selection, we can pick the best models during training via online evaluations (similar to early stopping in supervised learning). This regime provides a further boost in the performance of Diffusion-QL as shown in Table [4](#).

Tasks	learning rate	$\eta$	max Q backup
halfcheetah-medium-v2	$3 \times 10^{-4}$	1.0	False
hopper-medium-v2	$3 \times 10^{-4}$	1.0	False
walker2d-medium-v2	$3 \times 10^{-4}$	1.0	False
halfcheetah-medium-replay-v2	$3 \times 10^{-4}$	1.0	False
hopper-medium-replay-v2	$3 \times 10^{-4}$	1.0	False
walker2d-medium-replay-v2	$3 \times 10^{-4}$	1.0	False
halfcheetah-medium-expert-v2	$3 \times 10^{-4}$	1.0	False
hopper-medium-expert-v2	$3 \times 10^{-4}$	1.0	False
walker2d-medium-expert-v2	$3 \times 10^{-4}$	1.0	False
antmaze-umaze-v0	$3 \times 10^{-4}$	0.5	False
antmaze-umaze-diverse-v0	$3 \times 10^{-4}$	2.0	True
antmaze-medium-play-v0	$1 \times 10^{-3}$	2.0	True
antmaze-medium-diverse-v0	$3 \times 10^{-4}$	3.0	True
antmaze-large-play-v0	$3 \times 10^{-4}$	4.5	True
antmaze-large-diverse-v0	$3 \times 10^{-4}$	3.5	True
pen-human-v1	$3 \times 10^{-5}$	0.15	False
pen-cloned-v1	$3 \times 10^{-5}$	0.1	False
kitchen-complete-v0	$3 \times 10^{-4}$	0.005	False
kitchen-partial-v0	$3 \times 10^{-4}$	0.005	False
kitchen-mixed-v0	$3 \times 10^{-4}$	0.005	False

Table 3: Hyperparameter settings of all selected tasks.

Table 4: Performance comparison with online model selection and offline model selection.

AntMaze Tasks	Diffusion-QL (Offline)	Diffusion-QL (Online)
halfcheetah-medium-v2	51.1 $\pm$ 0.5	<b>51.5</b> $\pm$ 0.3
hopper-medium-v2	90.5 $\pm$ 4.6	<b>96.6</b> $\pm$ 3.4
walker2d-medium-v2	87.0 $\pm$ 0.9	<b>87.3</b> $\pm$ 0.5
halfcheetah-medium-replay-v2	47.8 $\pm$ 0.3	<b>48.3</b> $\pm$ 0.2
hopper-medium-replay-v2	101.3 $\pm$ 0.6	<b>102.0</b> $\pm$ 0.4
walker2d-medium-replay-v2	95.5 $\pm$ 1.5	<b>98.0</b> $\pm$ 0.5
halfcheetah-medium-expert-v2	96.8 $\pm$ 0.3	<b>97.2</b> $\pm$ 0.4
hopper-medium-expert-v2	111.1 $\pm$ 1.3	<b>112.3</b> $\pm$ 0.8
walker2d-medium-expert-v2	110.1 $\pm$ 0.3	<b>111.2</b> $\pm$ 0.9
<b>Average</b>	88.0	<b>89.3</b>
AntMaze Tasks	Diffusion-QL (Offline)	Diffusion-QL (Online)
antmaze-umaze-v0	93.4 $\pm$ 3.4	<b>96.0</b> $\pm$ 3.3
antmaze-umaze-diverse-v0	66.2 $\pm$ 8.6	<b>84.0</b> $\pm$ 10.1
antmaze-medium-play-v0	76.6 $\pm$ 10.8	<b>79.8</b> $\pm$ 8.7
antmaze-medium-diverse-v0	78.6 $\pm$ 10.3	<b>82.0</b> $\pm$ 9.5
antmaze-large-play-v0	46.4 $\pm$ 8.3	<b>49.0</b> $\pm$ 9.4
antmaze-large-diverse-v0	56.6 $\pm$ 7.6	<b>61.7</b> $\pm$ 8.2
<b>Average</b>	69.6	<b>75.4</b>
Adroit Tasks	Diffusion-QL (Offline)	Diffusion-QL (Online)
pen-human-v1	72.8 $\pm$ 9.6	<b>75.7</b> $\pm$ 9.0
pen-cloned-v1	57.3 $\pm$ 11.9	<b>60.8</b> $\pm$ 11.8
<b>Average</b>	65.1	<b>68.3</b>
Kitchen Tasks	Diffusion-QL (Offline)	Diffusion-QL (Online)
kitchen-complete-v0	84.0 $\pm$ 7.4	<b>84.5</b> $\pm$ 6.1
kitchen-partial-v0	60.5 $\pm$ 6.9	<b>63.7</b> $\pm$ 5.2
kitchen-mixed-v0	62.6 $\pm$ 5.1	<b>66.6</b> $\pm$ 3.3
<b>Average</b>	69.0	<b>71.6</b>