

**Table of Contents**

489	<b>List of Tables</b>	<b>1</b>
490	<b>List of Figures</b>	<b>1</b>
491	<b>A Environment &amp; Task Details</b>	<b>2</b>
492	A.1 MAZE . . . . .	2
493	A.2 FETCHPICK & FETCHPUSH . . . . .	2
494	A.3 HANDROTATE . . . . .	2
495	A.4 WALKER . . . . .	2
496	<b>B Algorithm</b>	<b>3</b>
497	<b>C Model Architecture</b>	<b>3</b>
498	C.1 Model Architecture of BC, Implicit BC, Diffusion Policy, and DBC . . . . .	3
499	C.2 Model Architecture of EBM, VAE, and GAN . . . . .	4
500	<b>D Training and Inference Details</b>	<b>5</b>
501	D.1 Computation Resource . . . . .	5
502	D.2 Hyperparameters . . . . .	5
503	D.3 Inference Details . . . . .	5
504	D.4 Comparing Different Generative Models . . . . .	6
505	<b>E Generalization Experiments in FETCHPICK and FETCHPUSH</b>	<b>8</b>
506	<b>F Qualitative Results and Additional Analysis</b>	<b>9</b>
507	F.1 Qualitative Results . . . . .	9
508	F.2 Learning Progress Analysis . . . . .	9
509	F.3 Episode Length Analysis of Goal-Directed Tasks . . . . .	10
510	<b>G On the Theoretical Motivation for Guiding Policy Learning with Diffusion Model</b>	<b>10</b>
511	<b>H Limitations</b>	<b>13</b>
512	<b>I Broader Impacts</b>	<b>13</b>

**List of Tables**

517	4	<b>Model Architectures.</b> . . . . .	4
518	5	<b>Hyperparameters.</b> . . . . .	5
519	6	<b>FETCHPICK Generalization Experimental Result.</b> . . . . .	8
520	7	<b>FETCHPUSH Generalization Experimental Result.</b> . . . . .	8
521	8	<b>Episode Length of Goal-Directed Tasks</b> . . . . .	11

**List of Figures**

523	5	<b>Qualitative Results</b> . . . . .	9
524	6	<b>Learning Progress.</b> . . . . .	10
525	7	<b>Episode Length of Goal-Directed Tasks</b> . . . . .	11
526	8	<b>Visualized Gradient Field</b> . . . . .	12

## 527 A Environment & Task Details

### 528 A.1 MAZE

529 **Description.** A point-maze agent in a 2D maze learns to navigate from its start location to a goal  
530 location by iteratively predicting its x and y acceleration. The 6D states include the agent’s two-  
531 dimensional current location and velocity, and the goal location. The start and the goal locations are  
532 randomized when an episode is initialized.

533 **Evaluation.** We evaluate the agents with 100 episodes and three random seeds and compare our  
534 method with the baselines regarding the average success rate and episode lengths, representing the  
535 effectiveness and efficiency of the policy learned by different methods. An episode terminates when  
536 the maximum episode length of 400 is reached.

537 **Expert Dataset.** The expert dataset consists of the 100 demonstrations with 18,525 transitions  
538 provided by Lee et al. [2021].

### 539 A.2 FETCHPICK & FETCHPUSH

540 **Description.** FETCHPICK requires a 7-DoF robot arm to pick up an object from the table and move  
541 it to a target location; FETCHPUSH requires the robot arm to push an object to a target location.  
542 Following the environment setups of Lee et al. [2021], a 16D state representation consists of the  
543 angles of the robot joints, the robot arm poses relative to the object, and goal locations. The first three  
544 dimensions of the action indicate the desired relative position at the next time step. For FETCHPICK,  
545 the fourth dimension of action specifies the distance between the two fingers of the gripper.

546 **Evaluation.** We evaluate the agents with 100 episodes and three random seeds and compare our  
547 method with the baselines regarding the average success rate and episode lengths. An episode  
548 terminates when the agent completes the task or the maximum episode length is reached, which is set  
549 to 50 for FETCHPICK and 120 for FETCHPUSH.

550 **Expert Dataset.** The expert dataset of FETCHPICK consists of 303 trajectories (10k transitions)  
551 while the expert dataset of FETCHPUSH consists of 185 trajectories (10k transitions) provided by Lee  
552 et al. [2021].

### 553 A.3 HANDROTATE

554 **Description.** HANDROTATE Plappert et al. [2018] requires a 24-DoF Shadow Dexterous Hand to  
555 in-hand rotate a block to a target orientation. The 68D state representation consists of the joint  
556 angles and velocities of the hand, object poses, and the target rotation. The 20D action indicates the  
557 position control of the 20 joints, which can be controlled independently. HANDROTATE is extremely  
558 challenging due to its high dimensional state and action spaces. We adapt the experimental setup  
559 used in Plappert et al. [2018] and Lee et al. [2021], where the rotation is restricted to the z-axis and  
560 the possible initial and target z rotations are set within  $[-\frac{\pi}{12}, \frac{\pi}{12}]$  and  $[\frac{\pi}{3}, \frac{2\pi}{3}]$ , respectively.

561 **Evaluation.** We evaluate the agents with 100 episodes and three random seeds and compare our  
562 method with the baselines regarding the average success rate and episode lengths. An episode  
563 terminates when the agent completes the goal or the maximum episode length of 50 is reached.

564 **Expert Dataset.** To collect expert demonstrations, we train a SAC Haarnoja et al. [2018] policy  
565 using dense rewards for 10M environment steps. The dense reward given at each time step  $t$  is  
566  $R(s_t, a_t) = d_t - d_{t+1}$ , where  $d_t$  and  $d_{t+1}$  represent the angles (in radian) between current and the  
567 desired block orientations before and after taking the actions. Following the training stage, the SAC  
568 expert policy achieves a success rate of 59.48%. Subsequently, we collect 515 successful trajectories  
569 (10k transitions) from this policy to form our expert dataset for HANDROTATE.

### 570 A.4 WALKER

571 **Description.** WALKER requires an agent to walk toward x-coordinate as fast as possible while  
572 maintaining its balance. The 17D state consists of angles of joints, angular velocities of joints, and  
573 velocities of the x and z-coordinate of the top. The 6D action specifies the torques to be applied on  
574 each joint of the walker avatar.

575 **Evaluation.** We evaluate each learned policy with 30 episodes and three random seeds and compare  
 576 our method with the baselines regarding the average returns of episodes and episode lengths. The  
 577 return of an episode is accumulated from all the time steps of an episode. An episode terminates when  
 578 the agent is unhealthy (*i.e.*, ill conditions predefined in the environment) or the maximum episode  
 579 length (1000) is reached.

580 **Expert Dataset.** The expert dataset consists of 5 trajectories with  $5k$  state-action pairs provided  
 581 by [Kostrikov \[2018\]](#).

## 582 B Algorithm

583 Our proposed framework DBC is detailed in Algorithm [1](#). The algorithm consists of two parts. (1)  
 584 **Learning a diffusion model:** The diffusion model  $\phi$  learns to model the distribution of concatenated  
 585 state-action pairs sampled from the demonstration dataset  $D$ . It learns to reverse the diffusion process  
 586 (*i.e.*, denoise) by optimizing  $\mathcal{L}_{\text{diff}}$ . (2) **Learning a policy with the learned diffusion model:** We  
 587 propose a diffusion model objective  $\mathcal{L}_{\text{DM}}$  for policy learning and jointly optimize it with the BC  
 588 objective  $\mathcal{L}_{\text{BC}}$ . Specifically,  $\mathcal{L}_{\text{DM}}$  is computed based on processing a sampled state-action pair  $(s, a)$   
 589 and a state-action pair  $(s, \hat{a})$  with the action  $\hat{a}$  predicted by the policy  $\pi$  with  $\mathcal{L}_{\text{diff}}$ .

---

### Algorithm 1 Diffusion Model-Augmented Behavioral Cloning (DBC)

---

**Input:** Expert’s Demonstration Dataset  $D$

**Output:** Policy  $\pi$ .

```

1: // Learning a diffusion model  $\phi$ 
2: Randomly initialize a diffusion model  $\phi$ 
3: for each diffusion model iteration do
4:   Sample  $(s, a)$  from  $D$ 
5:   Sample noise level  $n$  from  $\{0, \dots, N\}$ 
6:   Update  $\phi$  using  $L_{\text{diff}}$  from Eq. 2
7: end for
8: // Learning a policy  $\pi$  with the learned diffusion model  $\phi$ 
9: Randomly initialize a policy  $\pi$ 
10: for each policy iteration do
11:   Sample  $(s, a)$  from  $D$ 
12:   Predict an action  $\hat{a}$  using  $\pi$  from  $s$ :  $\hat{a} \sim \pi(s)$ 
13:   Compute the BC loss  $L_{\text{BC}}$  using Eq. 1
14:   Sample noise level  $n$  from  $\{0, \dots, N\}$ 
15:   Compute the agent diffusion loss  $L_{\text{diff}}^{\text{agent}}$  with  $(s, \hat{a})$  using Eq. 3
16:   Compute the expert diffusion loss  $L_{\text{diff}}^{\text{expert}}$  with  $(s, a)$  using Eq. 4
17:   Compute the diffusion model loss  $L_{\text{DM}}$  using Eq. 5
18:   Update  $\pi$  using the total loss  $L_{\text{total}}$  from Eq. 6
19: end for
20: return  $\pi$ 

```

---

## 590 C Model Architecture

591 This section describes the model architectures used for all the experiments. Section [C.1](#) presents the  
 592 model architectures of BC, Implicit BC, Diffusion Policy, and our proposed framework DBC. Section  
 593 [C.2](#) details the model architectures of the EBM, VAE, and GAN used for the experiment comparing  
 594 different generative models.

### 595 C.1 Model Architecture of BC, Implicit BC, Diffusion Policy, and DBC

596 We compare our DBC with three baselines (BC, Implicit BC, and Diffusion Policy) on various tasks  
 597 in Section [5.3](#). We detail the model architectures for all the methods on all the tasks in Table [4](#).  
 598 Note that all the models, the policy of BC, the energy-based model of Implicit BC, the conditional  
 599 diffusion model of Diffusion Policy, the policy and the diffusion model of DBC, are parameterized  
 600 by a multilayer perceptron (MLP). We report the implementation details for each method as follows.

Table 4: **Model Architectures.** We report the architectures used for all the methods on all the tasks.

Method	Models	Component	MAZE	FETHPICK	FETHPUSH	HANDROTATE	WALKER
BC	Policy $\pi$	# Layers	3	2	2	3	3
		Input Dim.	6	16	16	68	17
		Hidden Dim.	256	1024	1024	1024	256
		Output Dim.	2	4	3	20	6
Implicit BC	Policy $\pi$	# Layers	2	2	2	2	2
		Input Dim.	8	20	19	88	23
		Hidden Dim.	1024	1024	1024	512	1024
		Output Dim.	1	1	1	1	1
Diffusion Policy	Policy $\pi$	# Layers	5	5	5	5	5
		Input Dim.	8	20	19	88	23
		Hidden Dim.	256	1200	1024	2100	1200
		Output Dim.	2	4	3	20	6
DBC	DM $\phi$	# Layers	5	5	5	5	5
		Input Dim.	8	20	19	88	23
		Hidden Dim.	128	1024	1024	2048	1024
		Output Dim.	8	20	19	88	23
	Policy $\pi$	# Layers	3	2	2	3	3
		Input Dim.	6	16	16	68	17
		Hidden Dim.	256	1024	1024	512	256
		Output Dim.	2	4	3	20	6

601 **BC.** The non-linear activation function is a hyperbolic tangent for all the BC policies. We experiment  
 602 with BC policies with more parameters, which tend to severely overfit to expert datasets, resulting in  
 603 worse performance.

604 **Implicit BC.** The non-linear activation function is ReLU for all energy-based models of Implicit BC.  
 605 We empirically find that Implicit BC prefers shallow architectures in our tasks, so we set the number  
 606 of layers to 2 for the energy-based models.

607 **Diffusion Policy.** The non-linear activation function is ReLU for all the policies of Diffusion Policy.  
 608 We empirically find that Diffusion Policy performs better with a deeper architecture. Therefore, we  
 609 set the number of layers to 5 for the policy. In most cases, we use a Diffusion Policy with more  
 610 parameters than the total parameters of DBC consisting of the policy and the diffusion model.

611 **DBC.** The non-linear activation function is ReLU for the diffusion models and is a hyperbolic tangent  
 612 for the policies. We apply batch normalization and dropout layers with a 0.2 ratio for the diffusion  
 613 models on FETHPICK and FETHPUSH.

## 614 C.2 Model Architecture of EBM, VAE, and GAN

615 We compare different generative models (*i.e.*, EBM, VAE, and GAN) on MAZE in Section 5.5 and  
 616 we report the model architectures used for the experiment in this section.

617 **Energy-Based Model.** An energy-based model (EBM) consists of 5 linear layers with ReLU  
 618 activation. The EBM takes a concatenated state-action pair with a dimension of 8 as input; the output  
 619 is a 1-dimensional vector representing the estimated energy values of the state-action pair. The size  
 620 of the hidden dimensions is 128.

621 **Variational Autoencoder.** The architecture of a variational autoencoder consists of an encoder  
 622 and a decoder. The inputs of the encoder are a concatenated state-action pair, and the outputs  
 623 are the predicted mean and variance, which parameterize a Gaussian distribution. We apply the  
 624 reparameterization trick [Kingma and Welling, 2014], sample features from the predicted Gaussian  
 625 distribution, and use the decoder to produce the reconstructed state-action pair. The encoder and the  
 626 decoder both consist of 5 linear layers with LeakyReLU [Xu et al., 2020] activation. The size of the  
 627 hidden dimensions is 128. That said, the encoder maps an 8-dimensional state-action pair to two  
 628 128-dimensional vectors (*i.e.*, mean and variance), and the decoder maps a sampled 128-dimensional  
 629 vector back to an 8-dimensional reconstructed state-action pair.

Table 5: **Hyperparameters.** This table reports the hyperparameters used for all the methods on all the tasks. Note that our proposed framework (DBC) consists of two learning modules, the diffusion model and the policy, and therefore their hyperparameters are reported separately.

$\lambda$	Hyperparameter	MAZE	FETCHPICK	FETCHPUSH	HANDROTATE	WALKER
BC	Learning Rate	1e-4	1e-5	1e-5	5e-6	1e-4
	Batch Size	128	128	128	128	128
	# Epochs	2000	5000	5000	5000	2000
Implicit BC	Learning Rate	1e-4	5e-6	1e-4	1e-5	1e-4
	Batch Size	128	512	512	512	128
	# Epochs	10000	15000	15000	5000	10000
Diffusion Policy	Learning Rate	2e-4	1e-5	1e-5	1e-4	1e-4
	Batch Size	128	128	128	128	128
	# Epochs	20000	15000	15000	30000	10000
DBC (Ours)	Diffusion Model Learning rate	1e-3	1e-4	1e-4	3e-5	2e-4
	Diffusion Model Batch Size	128	128	128	128	1024
	Diffusion Model # Epochs	8000	10000	10000	10000	8000
	Policy Learning Rate	1e-4	1e-5	2e-5	1e-4	1e-4
	Policy Batch Size	128	128	128	128	128
	Policy # Epochs	2000	5000	5000	5000	2000
	$\lambda$	5	0.1	0.2	1	0.05

630 **Generative Adversarial Network.** The architecture of the generative adversarial network consists of  
631 a generator and a discriminator. The generator is the policy model that predicts an action from a given  
632 state, whose input dimension is 6 and output dimension is 2. On the other hand, the discriminator  
633 learns to distinguish the expert state-action pairs  $(s, a)$  from the state-action pairs produced by the  
634 generator  $(s, \hat{a})$ . Therefore, the input dimension of the discriminator is 8, and the output is a scalar  
635 representing the probability of the state-action pair being "real." The generator and the discriminator  
636 both consist of three linear layers with ReLU activation, and the size of the hidden dimensions is 256.

## 637 D Training and Inference Details

638 We describe the details of training and performing inference in this section, including computation  
639 resources and hyperparameters.

### 640 D.1 Computation Resource

641 We conducted all the experiments on the following three workstations:

- 642 • M1: ASUS WS880T workstation with an Intel Xeon W-2255 (10C/20T, 19.25M, 4.5GHz) 48-  
643 Lane CPU, 64GB memory, an NVIDIA RTX 3080 Ti GPU, and an NVIDIA RTX 3090 Ti  
644 GPU
- 645 • M2: ASUS WS880T workstation with an Intel Xeon W-2255 (10C/20T, 19.25M, 4.5GHz) 48-  
646 Lane CPU, 64GB memory, an NVIDIA RTX 3080 Ti GPU, and an NVIDIA RTX 3090 Ti  
647 GPU
- 648 • M3: ASUS WS880T workstation with an Intel Xeon W-2255 (10C/20T, 19.25M, 4.5GHz)  
649 48-Lane CPU, 64GB memory, and two NVIDIA RTX 3080 Ti GPUs

### 650 D.2 Hyperparameters

651 We report the hyperparameters used for all the methods on all the tasks in Table 5. We use the Adam  
652 optimizer [Kingma and Ba, 2015] for all the methods on all the tasks and use linear learning rate  
653 decay for all policy models.

### 654 D.3 Inference Details

655 This section describes how each method infers an action  $\hat{a}$  given a state  $s$ .

656 **BC & DBC.** The policy models of BC and DBC can directly predict an action given a state, *i.e.*,  
 657  $\hat{a} \sim \pi(s)$ , and are therefore more efficient during inference as described in Section 5.3.

658 **Implicit BC.** The energy-based model (EBM) of Implicit BC learns to predict an estimated energy  
 659 value for a state-action pair during training. To generate a predicted  $\hat{a}$  given a state  $s$  during inference,  
 660 it requires a procedure to sample and optimize actions. We follow Florence et al. [2022] and  
 661 implement a derivative-free optimization algorithm to perform inference.

662 The algorithm first randomly samples  $N_s$  vectors from the action space as candidates. The EBM then  
 663 produces the estimated energy value of each candidate action and applies the Softmax function on  
 664 the estimated energy values to produce a  $N_s$ -dimensional probability. Then, it samples candidate  
 665 actions according to the above probability and adds noise to them to generate another  $N_s$  candidates  
 666 for the next iteration. The above procedure iterates  $N_{iter}$  times. Finally, the action with maximum  
 667 probability in the last iteration is selected as the predicted action  $\hat{a}$ . In our experiments,  $N_s$  is set to  
 668 1000 and  $N_{iter}$  is set to 3.

669 **Diffusion Policy.** Diffusion Policy learns a conditional diffusion model as a policy and produces an  
 670 action from sampled noise vectors conditioning on the given state during inference. We follow Pearce  
 671 et al. [2023], Chi et al. [2023] and adopt Denoising Diffusion Probabilistic Models (DDPMs) J Ho  
 672 [2020] for the diffusion models. Once learned, the diffusion policy  $\pi$  can "denoise" a noise sampled  
 673 from a Gaussian distribution  $\mathcal{N}(0, 1)$  given a state  $s$  and yield a predicted action  $\hat{a}$  using the following  
 674 equation:

$$a_{n-1} = \frac{1}{\sqrt{\alpha_n}} \left( a_n - \frac{1 - \alpha_n}{\sqrt{1 - \alpha_n}} \pi(s, a_n, n) \right) + \sigma_n z, \quad (7)$$

675 where  $\alpha_n$ ,  $\bar{\alpha}_n$ , and  $\sigma_n$  are schedule parameters,  $n$  is the current time step of the reverse diffusion  
 676 process, and  $z \sim \mathcal{N}(0, 1)$  is a random vector. The above denoising process iterates  $N$  times to  
 677 produce a predicted action  $a_0$  from a sampled noise  $a_N \sim \mathcal{N}(0, 1)$ . The number of total diffusion  
 678 steps  $N$  is 100 in our experiment, which is the same for the diffusion model in DBC.

## 679 D.4 Comparing Different Generative Models

680 Our proposed framework employs a diffusion model (DM) to model the joint probability of expert  
 681 state-action pairs and utilizes it to guide policy learning. To justify our choice of generative models, we  
 682 explore using other popular generative models to replace the diffusion model in MAZE. Specifically,  
 683 we consider energy-based models (EBMs) [Du and Mordatch, 2019, Song and Kingma, 2021],  
 684 variational autoencoders (VAEs) [Kingma and Welling, 2014], and generative adversarial networks  
 685 (GANs) [Goodfellow et al., 2014]. Each generative model learns to model the joint distribution of  
 686 expert state-action pairs. For fair comparisons, all the policy models learning from learned generative  
 687 models consists of 3 linear layers with ReLU activation, where the hidden dimension is 256. All the  
 688 policies are trained for 2000 epochs using the Adam optimizer [Kingma and Ba, 2015], and a linear  
 689 learning rate decay is applied for EBMs and VAEs.

### 690 D.4.1 Energy-Based Model

691 **Model Learning.** Energy-based models (EBMs) learn to model the joint distribution of the expert  
 692 state-action pairs by predicting an estimated energy value for a state-action pair  $(s, a)$ . The EBM  
 693 aims to assign low energy value to the real expert state-action pairs while high energy otherwise.  
 694 Therefore, the predicted energy value can be used to evaluate how well a state-action pair  $(s, a)$  fits  
 695 the distribution of the expert state-action pair distribution.

696 To train the EBM, we generate  $N_{neg}$  random actions as negative samples for each expert state-action  
 697 pair as proposed in Florence et al. [2022]. The objective of the EBM  $E_\phi$  is the InfoNCE loss Oord  
 698 et al. [2018]:

$$\mathcal{L}_{\text{InfoNCE}} = \frac{e^{-E_\phi(s,a)}}{e^{-E_\phi(s,a)} + \sum_{i=1}^{N_{neg}} e^{-E_\phi(s,\tilde{a}_i)}}, \quad (8)$$

699 where  $(s, a)$  indicates an expert state-action pair,  $\tilde{a}_i$  indicates the sampled random action, and  $N_{neg}$   
 700 is set to 64 in our experiments. The EBM learns to separate the expert state-action pairs from the  
 701 negative samples by optimizing the above InfoNCE loss.

702 The EBM is trained for 8000 epochs with the Adam optimizer [Kingma and Ba, 2015], with a batch  
 703 size of 128 and an initial learning rate of 0.0005. We apply learning rate decay by 0.99 for every 100  
 704 epoch.

705 **Guiding Policy Learning.** To guide a policy  $\pi$  to learn, we design an EBM loss  $\mathcal{L}_{\text{EBM}} = E_{\phi}(s, \hat{a})$ ,  
 706 where  $\hat{a}$  indicates the predicted action produced by the policy. The above EBM loss regularizes the  
 707 policy to generate actions with low energy values, which encourage the predicted state-action pair  
 708  $(s, \hat{a})$  to fit the modeled expert state-action pair distribution. The policy learning from this EBM loss  
 709  $\mathcal{L}_{\text{EBM}}$  achieves a success rate of 49.09% in MAZE as reported in Table 2.

710 We also experiment with combining this EBM loss  $\mathcal{L}_{\text{EBM}}$  with the  $\mathcal{L}_{\text{BC}}$  loss. The policy optimizes  
 711  $\mathcal{L}_{\text{BC}} + \lambda_{\text{EBM}}\mathcal{L}_{\text{EBM}}$ , where  $\lambda_{\text{EBM}}$  is set to 0.1. Optimizing this combined loss yields a success rate of  
 712 80.00% in MAZE as reported in Table 2.

#### 713 D.4.2 Variational Autoencoder

714 **Model Learning.** Variational autoencoders (VAEs) model the joint distribution of the expert data  
 715 by learning to reconstruct expert state-action pairs  $(s, a)$ . Once the VAE is learned, how well a  
 716 state-action pair fits the expert distribution can be reflected in the reconstruction loss.

717 The objective of training a VAE is as follows:

$$\mathcal{L}_{\text{vae}} = \|\hat{x} - x\|^2 + D_{\text{KL}}(\mathcal{N}(\mu_x, \sigma_x) \|\mathcal{N}(0, 1)), \quad (9)$$

718 where  $x$  is the latent variable, *i.e.*, the concatenated state-action pair  $x = [s, a]$ , and  $\hat{x}$  is the  
 719 reconstruction of  $x$ , *i.e.*, the reconstructed state-action pair. The first term is the reconstruction loss,  
 720 while the second term encourages aligning the data distribution with a normal distribution  $\mathcal{N}(0, 1)$ ,  
 721 where  $\mu_x$  and  $\sigma_x$  are the predicted mean and standard deviation given  $x$ .

722 The VAE is trained for  $100k$  update iterations with the Adam optimizer [Kingma and Ba, 2015], with  
 723 a batch size of 128 and an initial learning rate of 0.0001. We apply learning rate decay by 0.5 for  
 724 every  $5k$  epoch.

725 **Guiding Policy Learning.** To guide a policy  $\pi$  to learn, we design a VAE loss  $\mathcal{L}_{\text{VAE}} = \max(\mathcal{L}_{\text{vae}}^{\text{agent}} -$   
 726  $\mathcal{L}_{\text{vae}}^{\text{expert}}, 0)$ , similar to Eq. 5. This loss forces the policy to predict an action, together with the state,  
 727 that can be well reconstructed with the learned VAE. The policy learning from this VAE loss  $\mathcal{L}_{\text{VAE}}$   
 728 achieves a success rate of 48.47% in MAZE as reported in Table 2.

729 We also experiment with combining this VAE loss  $\mathcal{L}_{\text{VAE}}$  with the  $\mathcal{L}_{\text{BC}}$  loss. The policy optimizes  
 730  $\mathcal{L}_{\text{BC}} + \lambda_{\text{VAE}}\mathcal{L}_{\text{VAE}}$ , where  $\lambda_{\text{VAE}}$  is set to 1. Optimizing this combined loss yields a success rate of  
 731 82.31% in MAZE as reported in Table 2.

#### 732 D.4.3 Generative Adversarial Network

733 **Adversarial Model Learning & Policy Learning.** Generative adversarial networks (GANs) model  
 734 the joint distribution of expert data with a generator and a discriminator. The generator aims to  
 735 synthesize a predicted action  $\hat{a}$  given a state  $s$ . On the other hand, the discriminator aims to identify  
 736 expert the state-action pair  $(s, a)$  from the predicted one  $(s, \hat{a})$ . Therefore, a learned discriminator  
 737 can evaluate how well a state-action pair fits the expert distribution.

738 While it is possible to learn a GAN separately and utilize the discriminator to guide policy learning,  
 739 we let the policy  $\pi$  be the generator directly and optimize the policy with the discriminator iteratively.  
 740 We hypothesize that a learned discriminator may be too selective for a policy training from scratch,  
 741 so we learn the policy  $\pi$  with the discriminator  $D$  to improve the policy and the discriminator  
 742 simultaneously.

743 The objective of training the discriminator  $D$  is as follows:

$$\mathcal{L}_{\text{disc}} = \text{BCE}(D(s, a), 1) + \text{BCE}(D(s, \hat{a}), 0) = -\log(D(s, a)) - \log(1 - D(s, \hat{a})), \quad (10)$$

744 where  $\hat{a} = \pi(s)$  is the predicted action, and  $\text{BCE}$  is the binary cross entropy loss. The binary label  
 745  $(0, 1)$  indicates whether or not the state-action pair sampled from the expert data. The generator and  
 746 the discriminator are both updated by Adam optimizers using a 0.00005 learning rate.

747 To learn a policy (*i.e.*, generator), we design the following GAN loss:

$$\mathcal{L}_{\text{GAN}} = \text{BCE}(D(s, \hat{a}), 1) = -\log(D(s, \hat{a})). \quad (11)$$

Table 6: **FETCHPICK Generalization Experimental Result.** We report the performance of our proposed framework DBC and the baselines regarding the mean and the standard deviation of the success rate with different levels of noise injected into the initial state and goal locations in FETCHPICK, evaluated over three random seeds.

Method	Noise Level				
	1	1.25	1.5	1.75	2
BC	86.78% $\pm$ 4.68%	69.15% $\pm$ 5.00%	54.42% $\pm$ 3.89%	43.49% $\pm$ 4.68%	36.64% $\pm$ 3.85%
Implicit BC	89.40% $\pm$ 4.85%	72.27% $\pm$ 6.71%	46.32% $\pm$ 5.49%	34.60% $\pm$ 4.78%	25.84% $\pm$ 4.16%
Diffusion Policy	76.04% $\pm$ 3.12%	74.37% $\pm$ 3.80%	69.22% $\pm$ 5.23%	56.95% $\pm$ 4.63%	53.93% $\pm$ 4.49%
DBC (Ours)	<b>97.59%</b> $\pm$ 1.53%	<b>88.71%</b> $\pm$ 6.46%	<b>78.76%</b> $\pm$ 10.84%	<b>69.36%</b> $\pm$ 12.72%	<b>62.62%</b> $\pm$ 14.01%

Table 7: **FETCHPUSH Generalization Experimental Result.** We report the performance of our proposed framework DBC and the baselines regarding the mean and the standard deviation of the success rate with different levels of noise injected into the initial state and goal locations in FETCHPUSH, evaluated over three random seeds.

Method	Noise Level				
	1	1.25	1.5	1.75	2
BC	94.07% $\pm$ 4.45%	82.52% $\pm$ 5.46%	66.02% $\pm$ 6.88%	48.85% $\pm$ 8.65%	34.82% $\pm$ 7.13%
Implicit BC	85.95% $\pm$ 8.39%	83.99% $\pm$ 6.06%	77.70% $\pm$ 4.42%	70.33% $\pm$ 6.06%	56.98% $\pm$ 11.74%
Diffusion Policy	97.92% $\pm$ 1.10%	93.02% $\pm$ 2.36%	86.93% $\pm$ 3.26%	74.50% $\pm$ 3.66%	65.84% $\pm$ 3.81%
DBC (Ours)	<b>99.83%</b> $\pm$ 0.23%	<b>99.38%</b> $\pm$ 0.78%	<b>94.92%</b> $\pm$ 3.09%	<b>87.48%</b> $\pm$ 5.04%	<b>78.43%</b> $\pm$ 7.41%

748 The above GAN loss guides the policy to generate state-action pairs that fit the joint distribution of  
 749 the expert data. The policy learning from this GAN loss  $\mathcal{L}_{\text{GAN}}$  achieves a success rate of 50.29% in  
 750 MAZE as reported in Table 2.

751 We also experiment with combining this GAN loss  $\mathcal{L}_{\text{GAN}}$  with the  $\mathcal{L}_{\text{BC}}$  loss. The policy optimizes  
 752  $\mathcal{L}_{\text{BC}} + \lambda_{\text{GAN}}\mathcal{L}_{\text{GAN}}$ , where  $\lambda_{\text{GAN}}$  is set to 0.2. Optimizing this combined loss yields a success rate of  
 753 71.64% in MAZE as reported in Table 2.

## 754 E Generalization Experiments in FETCHPICK and FETCHPUSH

755 This section further investigates the generalization capabilities of the policies learned by our proposed  
 756 framework and the baselines. To this end, we evaluate the policies by injecting different noise levels to  
 757 both the initial state and goal location in FETCHPICK and FETCHPUSH. Specifically, we parameterize  
 758 the noise by scaling the 2D sampling regions for the block and goal locations in both environments.  
 759 We expect all the methods to perform worse with higher noise levels, while the performance drop of  
 760 the methods with better generalization ability is less significant. The results are presented in Table 6  
 761 for FETCHPICK and Table 7 for FETCHPUSH.

762 **Overall Performance.** Our proposed framework DBC consistently outperforms all the baselines  
 763 with different noise levels, indicating the superiority of DBC when different levels of generalization  
 764 are required.

765 **Performance Drop with Increased Noise Level.** In FETCHPICK, DBC experiences a performance  
 766 drop of 35.0% when the noise level increase from 1 to 2. However, BC and Implicit BC demonstrate  
 767 a more significant performance drop of 50.1% and 63.6%, respectively. Notably, Diffusion Policy  
 768 initially performs poorly at a noise level of 1 but demonstrates its robustness with a performance  
 769 drop of only 22.1% when the noise level increases to 2. On the other hand, in FETCHPUSH, DBC  
 770 experiences a performance drop of 21.4% when the noise level increase from 1 to 2, while all the  
 771 baselines have a more significant performance drop: BC (63%), Implicit BC (33.7%), and Diffusion  
 772 Policy (32.8%). This demonstrates that our proposed framework not only generalizes better but also  
 773 exhibits greater robustness to noise compared to the baselines.

774 In this experiment, we set the coefficient  $\lambda$  of DBC to 0.5 in FETCHPUSH and 0.1 in FETCHPICK,  
 775 resulting in improved performance compared to the performance reported in the main paper.

776 In this experiment, we set the coefficient  $\lambda$  of DBC to 0.5 in FETCHPUSH and 0.1 in FETCHPICK,  
 777 resulting in an improved performance compared to the performance reported in the main paper.



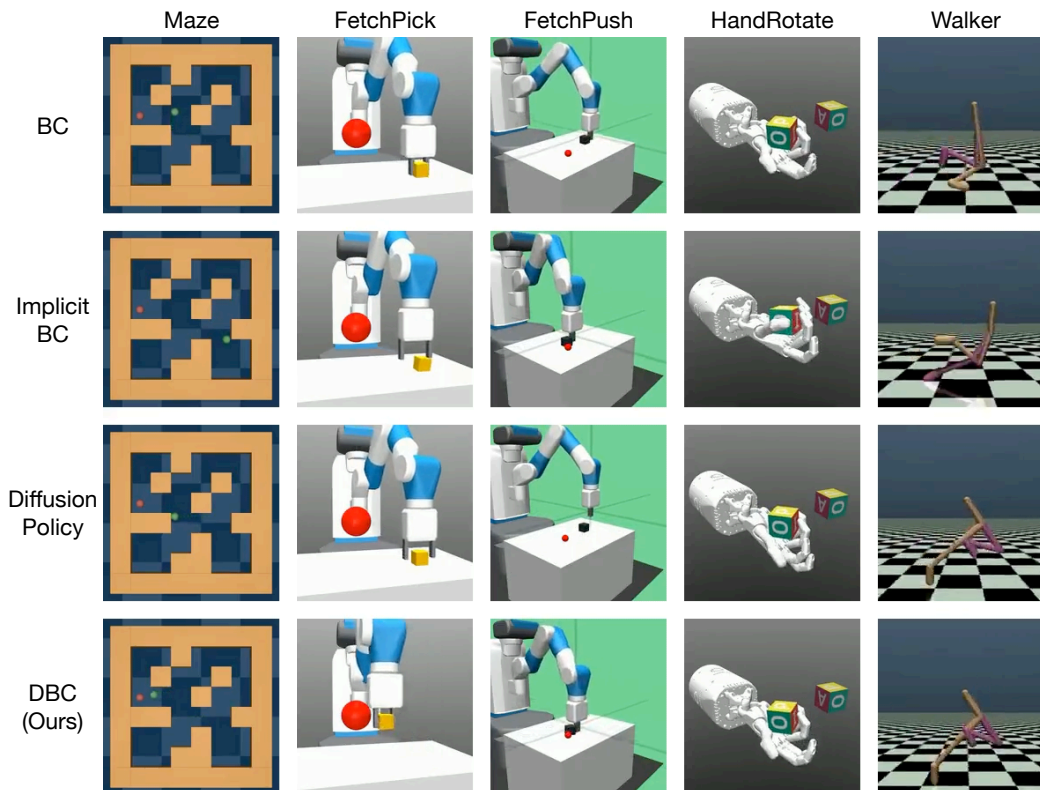


Figure 5: **Qualitative Results.** Rendered videos of the policies learned by our proposed framework and the baselines can be found at <https://sites.google.com/view/diffusion-behavioral-cloning>.

## 778 F Qualitative Results and Additional Analysis

779 This section provides more detailed analyses of our proposed framework and the baselines. We  
 780 present the qualitative results in Section F.1. Then, we analyze the learning progress and the episode  
 781 length of goal-directed tasks during inference in Section F.2 and Section F.3, respectively.

### 782 F.1 Qualitative Results

783 Rendered videos of the policies learned by our proposed framework and the baselines can be found  
 784 at <https://sites.google.com/view/diffusion-behavioral-cloning>. A screenshot of the  
 785 rendered videos on the web page is presented in Figure 5.

### 786 F.2 Learning Progress Analysis

787 In this section, we analyze the learning progress of all the methods on all the tasks. The training curves  
 788 are presented in Figure 6. Our proposed framework (DBC) not only achieves the best converged  
 789 performance but also converges the fastest, demonstrating its learning efficiency.

790 Since Implicit BC and Diffusion Policy take significantly longer to converge, we set a higher number  
 791 of training epochs for these two methods (see Table 5), and hence their learning curves are notably  
 792 longer than BC and DBC.

793 Note that we make sure the numbers of training epochs for Implicit BC and Diffusion Policy are not  
 794 less the total number of training epochs for learning both the diffusion model and the policy in DBC,  
 795 except for Implicit BC in HANDROTATE where training longer does not yield any improvement. This  
 796 forecloses the possibility of the superior performance of DBC coming from learning with a higher  
 797 total number of training epochs.

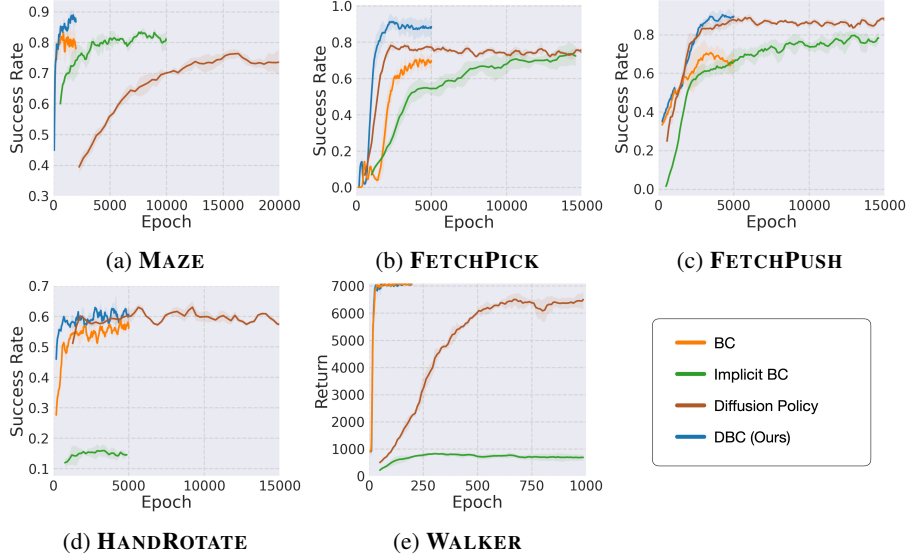


Figure 6: **Learning Progress.** We evaluate the baselines and our proposed method DBC and its variants during the learning process. Since Implicit BC (green) and Diffusion Policy (brown) take significantly longer to converge, we set a higher number of training epochs for these two methods, and hence their learning curves are notably longer than BC (orange) and DBC (blue). Our method demonstrates superior learning efficiency over the baselines.

### 798 F.3 Episode Length Analysis of Goal-Directed Tasks

799 In this section, we investigate the efficiency of the learned policies regarding the number of time steps  
800 they need to fulfill a task. We compare all the methods regarding average episode lengths over 100  
801 episodes and three random seeds in all goal-directed tasks (MAZE, FETCHPUSH, FETCHPICK, and  
802 HANDROTATE). The results are presented in Table 8 and Figure 7.

803 Note that Implicit BC and Diffusion Policy take significantly longer to converge, and hence we set  
804 a higher number of training epochs for these two methods (see Table 5). As a result, their learning  
805 curves are notably longer than BC and DBC.

806 We observe that our proposed framework DBC results in the shortest episode lengths in MAZE,  
807 FETCHPUSH, and FETCHPICK while performing competitively against the best-performing baseline  
808 (Diffusion Policy) in HANDROTATE. This indicates that DBC learns an efficient policy that can  
809 accomplish tasks quickly.

## 810 G On the Theoretical Motivation for Guiding Policy Learning with Diffusion 811 Model

812 This section further elaborates on the technical motivation for leveraging diffusion models for  
813 imitation learning. Specifically, we aim to learn a diffusion model to model the joint distribution of  
814 expert state-action pairs. Then, we propose to utilize this learned diffusion model to augment a BC  
815 policy that aims to imitate expert behaviors.

816 We consider the distribution of expert state-action pairs as the real data distribution  $q_x$  in learning a  
817 diffusion model. Following this setup,  $x_0$  represents an original expert state-action pair  $(s, a)$  and  
818  $q(x_n|x_{n-1})$  represents the forward diffusion process, which gradually adds Gaussian noise to the  
819 data in each timestep  $n = 1, \dots, N$  until  $x_N$  becomes an isotropic gaussian distribution. On the other  
820 hand, the reverse diffusion process is defined as  $\phi(x_{n-1}|x_n) := \mathcal{N}(x_{n-1}; \mu_\theta(x_n, n), \Sigma_\theta(x_n, n))$ ,  
821 where  $\theta$  denotes the learnable parameters of the diffusion model  $\phi$ , as illustrated in Figure 1.

822 Our key idea is to use the proposed diffusion model loss  $\mathcal{L}_{DM}$  in Eq. 5 as an estimate of how well a  
823 predicted state-action pair  $(s, \hat{a})$  fits the expert state-action pair distribution, as described in Section  
824 4.2.2. In the following derivation, we will show that by optimizing this diffusion model loss  $\mathcal{L}_{DM}$ , we

Table 8: **Episode Length of Goal-Directed Tasks.** We report the mean and the standard deviation of the episode length ( $\downarrow$ ) on MAZE, FETCHPICK, FETCHPUSH, and HANDROTATE, evaluated over three random seeds. The experiments demonstrate that our proposed method (DBC) outperforms (*i.e.*, finish tasks with fewer time steps) the baselines on MAZE, FETCHPICK, and FETCHPUSH while performing competitively in HANDROTATE.

Method	MAZE	FETCHPICK	FETCHPUSH	HANDROTATE
BC	219.95 $\pm$ 13.21	39.92 $\pm$ 0.65	74.08 $\pm$ 5.55	33.79 $\pm$ 1.18
Implicit BC	199.91 $\pm$ 15.95	44.67 $\pm$ 0.65	67.75 $\pm$ 3.13	46.13 $\pm$ 0.84
Diffusion Policy	241.45 $\pm$ 12.47	42.20 $\pm$ 0.64	80.93 $\pm$ 8.88	<b>31.95</b> $\pm$ 0.82
DBC (Ours)	<b>193.12</b> $\pm$ 10.30	<b>30.22</b> $\pm$ 1.38	<b>54.58</b> $\pm$ 3.33	<b>31.97</b> $\pm$ 1.49

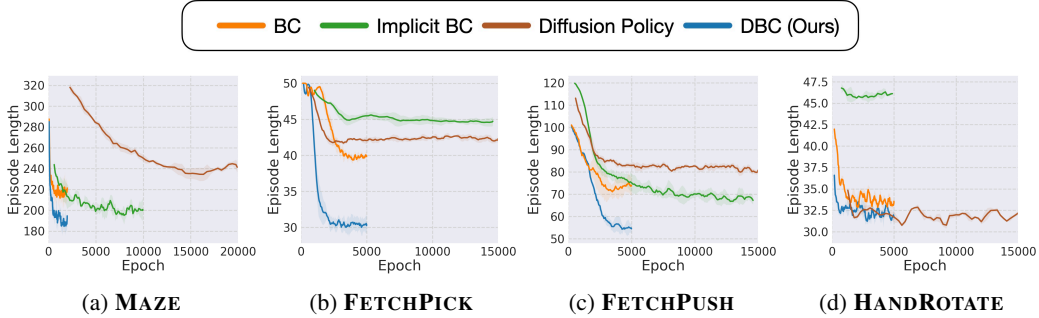


Figure 7: **Episode Length of Goal-Directed Tasks.** We evaluate the baselines and our proposed method regarding the episode length during the learning process. Since Implicit BC (green) and Diffusion Policy (brown) take significantly longer to converge, we set a higher number of training epochs for these two methods, and hence their learning curves are notably longer than BC (orange) and DBC (blue). The average episode length indicates how fast the agent reaches the goal, which can be a measurement of the efficiency of the agent. Our method DBC demonstrates superior efficiency in accomplishing tasks.

825 maximize the lower bound of the agent data’s probability under the derived expert distribution and  
 826 hence bring the agent policy  $\pi$  closer to the expert policy  $\pi^E$ , which is the goal of imitation learning.

827 As depicted in Luo [2022], one can conceptualize diffusion models, including DDPM [J Ho, 2020]  
 828 adopted in this work, as a hierarchical variational autoencoder [Kingma and Welling, 2014], which  
 829 maximizes the likelihood  $p(x)$  of observed data points  $x$ . Therefore, similar to hierarchical variational  
 830 autoencoders, diffusion models can optimize the Evidence Lower Bound (ELBO) by minimizing the  
 831 KL divergence  $D_{KL}(q(x_{n-1}|x_n, x_0) || \phi(x_{n-1}|x_n))$ . Consequently, this can be viewed as minimizing  
 832 the KL divergence to fit the distribution of the predicted state-action pairs  $(s, \hat{a})$  to the distribution of  
 833 expert state-action pairs.

834 According to Bayes’ theorem and the properties of Markov chains, the forward diffusion process  
 835  $q(x_{n-1}|x_n, x_0)$  follows:

$$q(x_{n-1}|x_n, x_0) \sim \mathcal{N}(x_{n-1}; \underbrace{\frac{\sqrt{\alpha_n}(1 - \bar{\alpha}_{n-1})x_n + \sqrt{\bar{\alpha}_{n-1}}(1 - \alpha_n)x_0}{1 - \bar{\alpha}_n}}_{\mu_q(x_n, x_0)}, \underbrace{\frac{(1 - \alpha_n)(1 - \bar{\alpha}_{n-1})}{1 - \bar{\alpha}_n}}_{\Sigma_q(n)}).$$

836 The variation term  $\Sigma_q(n)$  in the above equation can be written as  $\sigma_q^2(n)I$ , where  $\sigma_q^2(n) =$   
 837  $\frac{(1 - \alpha_n)(1 - \bar{\alpha}_{n-1})}{1 - \bar{\alpha}_n}$ . Therefore, minimizing the KL divergence is equivalent to minimizing the gap

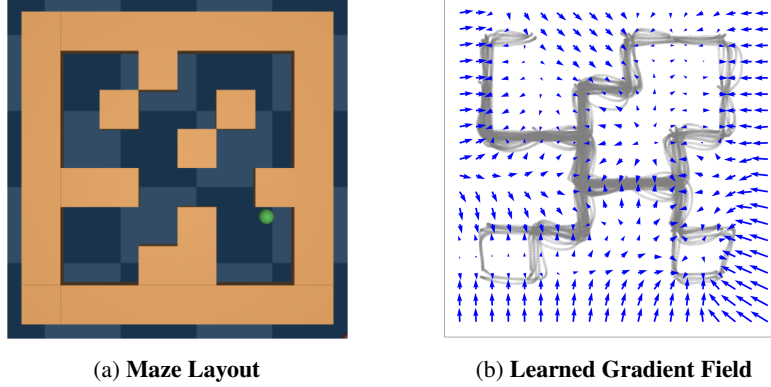


Figure 8: **Visualized Gradient Field.** (a) **Maze Layout:** The layout of the medium maze used for MAZE. (b) **Learned Gradient Field:** We visualize the MAZE expert demonstration as a distribution of points by their first two dimensions in gray. The points that cluster densely have a high probability, and vice versa. Once a diffusion model is well-trained, it can move randomly sampled points to the area with high probability by predicting gradients (blue arrows). Accordingly, the estimate  $p(s, a)$  of joint distribution modeling can serve as guidance for policy learning, as proposed in this work.

838 between the mean values of the two distributions:

$$\begin{aligned}
 & \arg \min_{\theta} D_{KL}(q(x_{n-1}|x_n, x_0) || \phi(x_{n-1}|x_n)) \\
 &= \arg \min_{\theta} D_{KL}(\mathcal{N}(x_{n-1}; \mu_q, \Sigma_q(n)) || \mathcal{N}(x_{n-1}; \mu_{\theta}, \Sigma_q(n))) \\
 &= \arg \min_{\theta} \frac{1}{2\sigma_q^2(n)} [||\mu_{\theta} - \mu_q||_2^2],
 \end{aligned}$$

839 where  $\mu_q$  represents the denoising transition mean and  $\mu_{\theta}$  represents the approximated denoising  
 840 transition mean by the model.

Different implementations adopt different forms to model  $\mu_{\theta}$ . Specifically, for DDPMs adopted in this work, the true denoising transition mean  $\mu_q(x_n, x_0)$  derived above can be rewritten as:

$$\mu_q(x_n, x_0) = \frac{1}{\sqrt{\alpha_n}} \left( x_n - \frac{1 - \alpha_n}{\sqrt{1 - \bar{\alpha}_n}} \epsilon_0 \right),$$

841 which is referenced from Eq. 11 in [J Ho \[2020\]](#). Hence, we can set our approximate denoising  
 842 transition mean  $\mu_{\theta}$  in the same form as the true denoising transition mean:

$$\mu_{\theta}(x_n, n) = \frac{1}{\sqrt{\alpha_n}} \left( x_n - \frac{1 - \alpha_n}{\sqrt{1 - \bar{\alpha}_n}} \hat{\epsilon}_{\theta}(x_n, n) \right), \tag{12}$$

843 as illustrated in [Popov et al. \[2022\]](#). [Song et al. \[2021\]](#) further show that the entire diffusion model  
 844 formulation can be revised to view continuous stochastic differential equations (SDEs) as a forward  
 845 diffusion. It points out that the reverse process is also an SDE, which can be computed by estimating  
 846 a score function  $\nabla_x \log p_t(x)$  at each denoising time step. The idea of representing a distribution  
 847 by modeling its score function is introduced in [Song and Ermon \[2019\]](#). The fundamental concept  
 848 is to model the gradient of the log probability density function  $\nabla_x \log p_t(x)$ , a quantity commonly  
 849 referred to as the (Stein) score function. Such score-based models are not required to have a tractable  
 850 normalizing constant and can be directly acquired through score matching. The measure of this score  
 851 function determines the optimal path to take in the space of the data distribution to maximize the log  
 852 probability under the derived real distribution.

853 As shown in Figure [8b](#), we visualized the learned gradient field of a diffusion model, which learns to  
 854 model the expert state-action pairs in MAZE. Once trained, this diffusion model can guide a policy  
 855 with predicted gradients (blue arrows) to move to areas with high probability, as proposed in our  
 856 work.

857 Essentially, by moving in the opposite direction of the source noise, which is added to a data point  $x_t$   
 858 to corrupt it, the data point is “denoised”; hence the log probability is maximized. This is supported

859 by the fact that modeling the score function is the same as modeling the negative of the source noise.  
 860 This perspective of the diffusion model is dubbed diffusion SDE. Moreover, [Popov et al. \(2022\)](#)  
 861 prove that Eq. [12](#) is diffusion SDE’s maximum likelihood SDE solver. Hence, the corresponding  
 862 divergence optimization problem can be rewritten as:

$$\begin{aligned} & \arg \min_{\theta} D_{KL}(q(x_{n-1}|x_n, x_0) || \phi(x_{n-1}|x_n)) \\ & = \arg \min_{\theta} \frac{1}{2\sigma_q^2(n)} \frac{(1 - \alpha_n)^2}{(1 - \bar{\alpha}_n)\alpha_n} [||\hat{\epsilon}_{\theta}(x_n, n) - \epsilon_0||_2^2], \end{aligned}$$

where  $\epsilon_{\theta}$  is a function approximator aim to predict  $\epsilon$  from  $x$ . As the coefficients can be omitted during optimization, we yield the learning objective  $\mathcal{L}_{\text{diff}}$  as stated in in Eq. [2](#):

$$\mathcal{L}_{\text{diff}} = ||\hat{\epsilon}(s, a, n) - \epsilon(n)||^2 = ||\phi(s, a, \epsilon(n)) - \epsilon(n)||^2.$$

863 The above derivation motivates our proposed framework that augments a BC policy by using the  
 864 diffusion model to provide guidance that captures the joint probability of expert state-action pairs.  
 865 Based on the above derivation, minimizing the proposed diffusion model loss (*i.e.*, learning to denoise)  
 866 is equivalent to finding the optimal path to take in the data space to maximize the log probability. To  
 867 be more accurate, when the learner policy predicts an action that obtains a lower  $\mathcal{L}_{\text{diff}}$ , it means that  
 868 the predicted action  $\hat{a}$ , together with the given state  $s$ , fits better with the expert distribution.

869 Accordingly, by minimizing our proposed diffusion loss, the policy is encouraged to imitate the  
 870 expert policy. To further alleviate the impact of rarely-seen state-action pairs  $(s, a)$ , we propose  
 871 to compute the above diffusion loss for both expert data  $(s, a)$  and predicted data  $(s, \hat{a})$  and yield  
 872  $\mathcal{L}_{\text{diff}}^{\text{expert}}$  and  $\mathcal{L}_{\text{diff}}^{\text{agent}}$ , respectively. Therefore, we propose to augment BC with this objective:  $\mathcal{L}_{\text{DM}} =$   
 873  $\max(\mathcal{L}_{\text{diff}}^{\text{agent}} - \mathcal{L}_{\text{diff}}^{\text{expert}}, 0)$  This design is justified in Section [5.6.2](#).

## 874 H Limitations

875 This section discusses the limitations of our proposed framework.

- 876 • Since this work aims to learn from demonstrations without interacting with environments, our  
 877 proposed framework in its current form is only designed to learn from expert trajectories and  
 878 cannot learn from trajectories produced by the learner policy. Extending our method to incorporate  
 879 agent data can potentially allow for improvement when interacting environments are possible,  
 880 which is left for future work.
- 881 • The key insight of our work is to allow the learner policy to benefit from both modeling the  
 882 conditional and joint probability of expert state-action distributions. To this end, we propose to  
 883 optimize both the BC loss and the proposed diffusion model loss. To balance the importance of  
 884 the two losses, we introduce a coefficient  $\lambda$  as an additional hyperparameter. While the ablation  
 885 study conducted in MAZE shows that the performance of our proposed framework is robust to  
 886  $\lambda$ , this can potentially increase the difficulty of searching for optimal hyperparameters when  
 887 applying our proposed framework to a new application.

## 888 I Broader Impacts

889 This work proposes Diffusion Model-Augmented Behavioral Cloning, a novel imitation learning  
 890 framework that aims to increase the ability of autonomous learning agents (*e.g.*, robots, game AI  
 891 agents) to acquire skills by imitating demonstrations provided by experts (*e.g.*, humans). However, it  
 892 is crucial to acknowledge that our proposed framework, by design, inherits any biases exhibited by  
 893 the expert demonstrators. These biases can manifest as sub-optimal, unsafe, or even discriminatory  
 894 behaviors. To address this concern, ongoing research endeavors to mitigate bias and promote  
 895 fairness in machine learning hold promise in alleviating these issues. Moreover, research works  
 896 that enhance learning agents’ ability to imitate experts, such as this work, can pose a threat to job  
 897 security. Nevertheless, in sum, we firmly believe that our proposed framework can offer tremendous  
 898 advantages in terms of enhancing the quality of human life and automating laborious, arduous, or  
 899 perilous tasks that pose risks to humans, which far outweigh the challenges and potential issues.

900