

---

# Diffusion Model-Augmented Behavioral Cloning

---

Anonymous Author(s)

Affiliation

Address

email

## Abstract

1 Imitation learning addresses the challenge of learning by observing an expert’s  
2 demonstrations without access to reward signals from environments. Most existing  
3 imitation learning methods that do not require interacting with environments either  
4 model the expert distribution as the conditional probability  $p(a|s)$  (e.g., behavioral  
5 cloning, BC) or the joint probability  $p(s, a)$  (e.g., implicit behavioral cloning). De-  
6 spite its simplicity, modeling the conditional probability with BC usually struggles  
7 with generalization. While modeling the joint probability can lead to improved  
8 generalization performance, the inference procedure can be time-consuming and it  
9 often suffers from manifold overfitting. This work proposes an imitation learning  
10 framework that benefits from modeling both the conditional and joint probability  
11 of the expert distribution. Our proposed diffusion model-augmented behavioral  
12 cloning (DBC) employs a diffusion model trained to model expert behaviors and  
13 learns a policy to optimize both the BC loss (conditional) and our proposed diffu-  
14 sion model loss (joint). DBC outperforms baselines in various continuous control  
15 tasks in navigation, robot arm manipulation, dexterous manipulation, and locomo-  
16 tion. We design additional experiments to verify the limitations of modeling either  
17 the conditional probability or the joint probability of the expert distribution as well  
18 as compare different generative models.

## 19 1 Introduction

20 Recently, the success of deep reinforcement learning (DRL) [Mnih et al., 2015, Lillicrap et al., 2016,  
21 Arulkumaran et al., 2017] has inspired the research community to develop DRL frameworks to  
22 control robots, aiming to automate the process of designing sensing, planning, and control algorithms  
23 by letting the robot learn in an end-to-end fashion. Yet, acquiring complex skills through trial and  
24 error can still lead to undesired behaviors even with sophisticated reward design [Christiano et al.,  
25 2017, Leike et al., 2018, Lee et al., 2019]. Moreover, the exploring process could damage expensive  
26 robotic platforms or even be dangerous to humans [Garcia and Fernández, 2015, Levine et al., 2020].

27 To overcome this issue, imitation learning (*i.e.*, learning from demonstration) [Schaal, 1997, Osa et al.,  
28 2018] has received growing attention, whose aim is to learn a policy from expert demonstrations,  
29 which are often more accessible than appropriate reward functions for reinforcement learning. Among  
30 various imitation learning directions, adversarial imitation learning [Ho and Ermon, 2016, Zolna  
31 et al., 2021, Kostrikov et al., 2019] and inverse reinforcement learning [Ng and Russell, 2000, Abbeel  
32 and Ng, 2004] have achieved encouraging results in a variety of domains. Yet, these methods require  
33 interacting with environments, which can still be expensive or unsafe.

34 On the other hand, behavioral cloning (BC) [Pomerleau, 1989, Bain and Sammut, 1995] does not  
35 require interacting with environments. BC formulates imitation learning as a supervised learning  
36 problem — given an expert demonstration dataset, an agent policy takes states sampled from the  
37 dataset as input and learns to replicate the corresponding expert actions. One can view a BC policy as

38 a discriminative model  $p(a|s)$  that models the *conditional probability* of an action  $a$  given a state  $s$ .  
39 Due to its simplicity and training stability, BC has been widely adopted for various applications.

40 However, BC struggles at generalizing to states unobserved during training [Nguyen et al., 2023].  
41 To address this issue, implicit behavioral cloning (IBC) [Florence et al., 2022] aims to model the  
42 *joint probability* of the expert state-action pairs  $p(s, a)$  with energy-based models. IBC demonstrates  
43 superior performance when generalization is required. Yet, imitation learning methods in a similar  
44 vein [Ganapathi et al., 2022] that model the *joint probability* of state-action pairs  $p(s, a)$  instead  
45 of directly predicting actions  $p(a|s)$  require time-consuming actions sampling and optimization to  
46 retrieve a desired action  $\arg \max_{a \in \mathcal{A}} p(s, a)$  during inference despite the choice of models.

47 This work proposes an imitation learning framework that combines both the efficiency of modeling the  
48 *conditional probability* and the generalization ability of modeling the *joint probability*. Specifically,  
49 we propose to model the expert state-action pairs using a state-of-the-art generative model, a diffusion  
50 model, which learns to estimate how likely a state-action pair is sampled from the expert dataset.  
51 Then, we train a policy to optimize both the BC objective and the estimate produced by the learned  
52 diffusion model. Therefore, our proposed framework not only can efficiently predict actions given  
53 states via capturing the *conditional probability*  $p(a|s)$  but also enjoys the generalization ability  
54 induced by modeling the *joint probability*  $p(s, a)$  and utilizing it to guide policy learning.

55 We evaluate our proposed framework and baselines in various continuous control domains, including  
56 navigation, robot arm manipulation, and locomotion. The experimental results show that the proposed  
57 framework outperforms all the baselines or achieves competitive performance on all tasks. Extensive  
58 ablation studies compare our proposed method to its variants, justifying our design choices, such as  
59 different generative models, and investigating the effect of hyperparameters.

## 60 2 Related Work

61 Imitation learning addresses the challenge of learning by observing expert demonstrations without  
62 access to reward signals from environments. It has various applications such as robotics [Schaal,  
63 1997], autonomous driving [Ly and Akhloofi, 2020], and game AI [Harmer et al., 2018].

64 **Behavioral Cloning (BC).** BC [Pomerleau, 1989, Torabi et al., 2018] formulate imitating an expert  
65 as a supervised learning problem. Due to its simplicity and effectiveness, it has been widely adopted  
66 in various domains. Yet, it often struggles at generalizing to states unobserved from the expert  
67 demonstrations [Ross et al., 2011, Florence et al., 2022]. In this work, we augment BC by employing  
68 a diffusion model that learns to capture the joint probability of expert state-action pairs.

69 **Adversarial Imitation Learning (AIL).** AIL methods aim to match the state-action distributions of  
70 an agent and an expert via adversarial training. Generative adversarial imitation learning (GAIL) [Ho  
71 and Ermon, 2016] and its extensions [Torabi et al., 2019, Kostrikov et al., 2019, Zolna et al., 2021]  
72 resemble the idea of generative adversarial networks [Goodfellow et al., 2014], which trains a  
73 generator policy to imitate expert behaviors and a discriminator to distinguish between the expert  
74 and the learner’s state-action pair distributions. While modeling state-action distributions often leads  
75 to satisfactory performance, adversarial learning can be unstable and inefficient [Chen et al., 2020].  
76 Moreover, AIL methods require online interaction with environments, which can be costly or even  
77 dangerous. In contrast, our work does not require interacting with environments.

78 **Inverse Reinforcement Learning (IRL).** IRL methods [Ng and Russell, 2000, Abbeel and Ng,  
79 2004, Fu et al., 2018, Lee et al., 2021] are designed to infer the reward function that underlies the  
80 expert demonstrations and then learn a policy using the inferred reward function. This allows for  
81 learning tasks whose reward functions are difficult to specify manually. However, due to its double-  
82 loop learning procedure, IRL methods are typically computationally expensive and time-consuming.  
83 Additionally, obtaining accurate estimates of the expert’s reward function can be difficult, especially  
84 when the expert’s behavior is non-deterministic or when the expert’s demonstrations are sub-optimal.

85 **Diffusion Policies.** Recently, Pearce et al. [2023], Chi et al. [2023], Reuss et al. [2023] propose to  
86 represent and learn an imitation learning policy using a conditional diffusion model, which produces  
87 a predicted action conditioning on a state and a sampled noise vector. These methods achieve  
88 encouraging results in modeling stochastic and multimodal behaviors from human experts or play

89 data. In contrast, instead of representing a policy using a diffusion model, our work employs a  
90 diffusion model trained on expert demonstrations to guide a policy as a learning objective.

## 91 3 Preliminaries

### 92 3.1 Imitation Learning

93 Without loss of generality, the reinforcement learning problem can be formulated as a Markov decision  
94 process (MDP), which can be represented by a tuple  $M = (S, A, R, P, \rho, \gamma)$  with states  $S$ , actions  
95  $A$ , reward function  $R(S, A) \in (0, 1)$ , transition distribution  $P(s' | s, a) : S \times A \times S \rightarrow [0, 1]$ , initial  
96 state distribution  $\rho$ , and discounted factor  $\gamma$ . Based on the rewards received while interacting with  
97 the environment, the goal is to learn a policy  $\pi(\cdot | s)$  to maximize the expectation of the cumulative  
98 discounted return (*i.e.*, value function):  $V(\pi) = \mathbb{E}[\sum_{t=0}^T \gamma^t R(s_t, a_t) | s_0 \sim \rho(\cdot), a_t \sim \pi(\cdot | s_t), s_{t+1} \sim$   
99  $P(s_{t+1} | s_t, a_t)]$ , where  $T$  denotes the episode length. Instead of interacting with the environment and  
100 receiving rewards, imitation learning aims to learn an agent policy from an expert demonstration  
101 dataset, containing  $M$  trajectories,  $D = \{\tau_1, \dots, \tau_M\}$ , where  $\tau_i$  represents a sequence of  $n_i$  state-  
102 action pairs  $\{s_1^i, a_1^i, \dots, s_{n_i}^i, a_{n_i}^i\}$ .

### 103 3.2 Behavioral Cloning: Modeling Conditional Probability $p(a|s)$

104 To learn a policy  $\pi$ , behavioral cloning (BC) directly estimates the expert policy  $\pi^E$  with maximum  
105 likelihood estimation (MLE). Given a state-action pair  $(s, a)$  sampled from the dataset  $D$ , BC  
106 optimizes  $\max_{\theta} \sum_{(s,a) \in D} \log(\pi_{\theta}(a|s))$ , where  $\theta$  denotes the parameters of the policy  $\pi$ . One can view a  
107 BC policy as a discriminative model  $p(a|s)$ , capturing the *conditional probability* of an action  $a$  given  
108 a state  $s$ . Despite its success in various applications, BC tends to overfit and struggle at generalizing  
109 to states unseen during training [Ross et al., 2011, Codevilla et al., 2019, Wang et al., 2022].

### 110 3.3 Modeling Joint Probability $p(s, a)$

111 Aiming for improved generalization ability, implicit behavioral cloning [Florence et al., 2022] and  
112 methods in a similar vein [Ganapathi et al., 2022] model the *joint probability*  $p(s, a)$  of expert state-  
113 action pairs. These methods demonstrate superior generalization performance in diverse domains. Yet,  
114 without directly modeling the *conditional probability*  $p(a|s)$ , the action sampling and optimization  
115 procedure to retrieve a desired action  $\arg \max_{a \in \mathcal{A}} p(s, a)$  during inference is often time-consuming.

116 Moreover, explicit generative models such as energy-based models [Du and Mordatch, 2019, Song and  
117 Kingma, 2021], variational autoencoder [Kingma and Welling, 2014], and flow-based models Rezende  
118 and Mohamed [2015], Dinh et al. [2017] are known to struggle with modeling observed high-  
119 dimensional data that lies on a low-dimensional manifold (*i.e.*, manifold overfitting) [Wu et al., 2021,  
120 Loaiza-Ganem et al., 2022]. As a result, these methods often perform poorly when learning from  
121 demonstrations produced by script policies or PID controllers, as discussed in Section 5.4.

122 We aim to develop an imitation learning framework that enjoys the advantages of modeling the  
123 *conditional probability*  $p(a|s)$  and the *joint probability*  $p(s, a)$ . Specifically, we propose to model the  
124 *joint probability* of expert state-action pairs using an explicit generative model  $\phi$ , which learns to  
125 produce an estimate indicating how likely a state-action pair is sampled from the expert dataset. Then,  
126 we train a policy to model the *conditional probability*  $p(a|s)$  by optimizing the BC objective and  
127 the estimate produced by the learned generative model  $\phi$ . Hence, our method can efficiently predict  
128 actions given states, generalize better to unseen states, and suffer less from manifold overfitting.

### 129 3.4 Diffusion Models

130 As described in the previous sections, this work aims to combine the advantages of modeling both  
131 the *conditional probability*  $p(a|s)$  and the *joint probability*  $p(s, a)$ . To this end, we leverage diffusion  
132 models to model the *joint probability* of expert state-action pairs. The diffusion model is a recently  
133 developed class of generative models and has achieved state-of-the-art performance on various  
134 tasks Sohl-Dickstein et al. [2015], Nichol and Dhariwal [2021], Dhariwal and Nichol [2021].

135 In this work, we utilize Denoising Diffusion  
 136 Probabilistic Models (DDPMs) J Ho [2020]  
 137 to model expert state-action pairs. Specifically,  
 138 DDPM models gradually add noise to data sam-  
 139 ples (*i.e.*, concatenated state-action pairs)  
 140 until they become isotropic Gaussian (*forward dif-*  
 141 *fusion process*), and then learn to denoise each  
 142 step and restore the original data samples (*re-*  
 143 *verse diffusion process*), as illustrated in Figure  
 144 1. In other words, DDPM learns to recognize a  
 145 data distribution by learning to denoise noisy  
 146 sampled data. More discussion on diffusion  
 147 models can be found in the Section G.

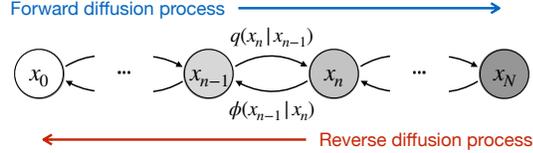


Figure 1: **Denoising Diffusion Probabilistic Model (DDPM).** Latent variables  $x_1, \dots, x_N$  are produced from the data point  $x_0$  via the forward diffusion process, *i.e.*, gradually adding noises to the latent variables. The diffusion model  $\phi$  learns to reverse the diffusion process by denoising the noisy data to reconstruct the original data point  $x_0$ .

## 148 4 Approach

149 Our goal is to design an imitation learning framework that enjoys both the advantages of modeling  
 150 the *conditional probability* and the *joint probability* of expert behaviors. To this end, we first adopt  
 151 behavioral cloning (BC) for modeling the *conditional probability* from expert state-action pairs, as  
 152 described in Section 4.1. To capture the *joint probability* of expert state-action pairs, we employ  
 153 a diffusion model which learns to produce an estimate indicating how likely a state-action pair is  
 154 sampled from the expert state-action pair distribution, as presented in Section 4.2.1. Then, we propose  
 155 to guide the policy learning by optimizing this estimate provided by a learned diffusion model,  
 156 encouraging the policy to produce actions similar to expert actions, as discussed in Section 4.2.2.  
 157 Finally, in Section 4.3, we introduce the framework that combines the BC loss and our proposed  
 158 diffusion model loss, allowing for learning a policy that benefits from modeling both the *conditional*  
 159 *probability* and the *joint probability* of expert behaviors. An overview of our proposed framework is  
 160 illustrated in Figure 2, and the algorithm is detailed in Section B.

### 161 4.1 Behavioral Cloning Loss

162 The behavioral cloning (BC) model aims to imitate expert behaviors with supervised learning. BC  
 163 learns to capture the conditional probability  $p(a|s)$  of expert state-action pairs. Given a sampled  
 164 expert state-action pair  $(s, a)$ , a policy  $\pi$  learns to predict an action  $\hat{a} \sim \pi(s)$  by optimizing

$$\mathcal{L}_{BC} = d(a, \hat{a}), \quad (1)$$

165 where  $d(\cdot, \cdot)$  denotes a distance measure between a pair of actions. For example, we can adapt the  
 166 mean-square error (MSE) loss  $\|a - \hat{a}\|^2$  for most continuous control tasks.

### 167 4.2 Learning a Diffusion Model and Guiding Policy Learning

168 Instead of directly learning the conditional probability  $p(a|s)$ , this section discusses how to model  
 169 the joint probability  $p(s, a)$  of expert behaviors with a diffusion model in Section 4.2.1 and presents  
 170 how to leverage the learned diffusion model to guide policy learning in Section 4.2.2.

#### 171 4.2.1 Learning a Diffusion Model

172 We propose to model the joint probability of expert state-action pairs with a diffusion model  $\phi$ .  
 173 Specifically, we create a joint distribution by simply concatenating a state vector  $s$  and an action  
 174 vector  $a$  from a state-action pair  $(s, a)$ . To model such distribution by learning a denoising diffusion  
 175 probabilistic model (DDPM) J Ho [2020], we inject noise  $\epsilon(n)$  into sampled state-action pairs, where  
 176  $n$  indicates the number of steps of the Markov procedure, which can be viewed as a variable of the  
 177 level of noise. Then, we train the diffusion model  $\phi$  to predict the injected noises by optimizing

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

178 where  $\hat{\epsilon}$  is the noise predicted by the diffusion model  $\phi$ . Once optimized, the diffusion model can  
 179 *recognize* the expert distribution by perfectly predicting the noise injected into state-action pairs  
 180 sampled from the expert distribution. On the other hand, predicting the noise injected into state-  
 181 action pairs sampled from any other distribution should yield a higher loss value. Therefore, we

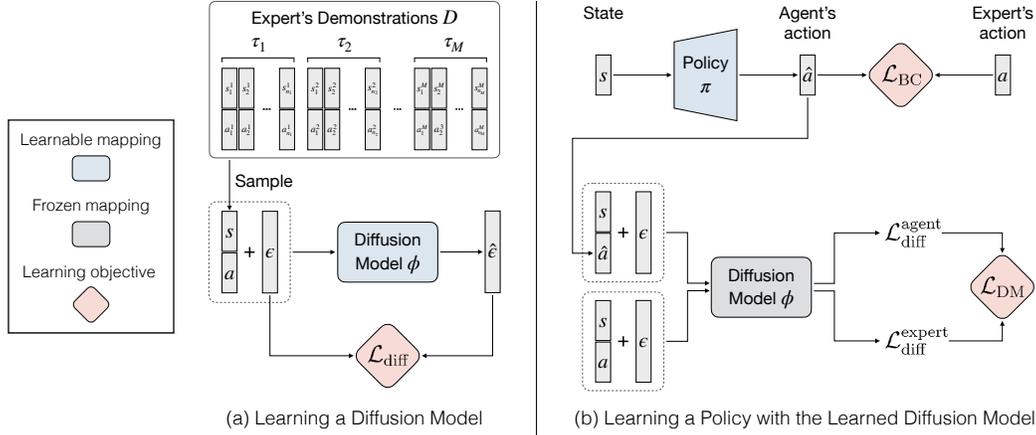


Figure 2: **Diffusion Model-Augmented Behavioral Cloning.** Our proposed method DBC augments behavioral cloning (BC) by employing a diffusion model. (a) **Learning a Diffusion Model:** the diffusion model  $\phi$  learns to model the distribution of concatenated state-action pairs sampled from the demonstration dataset  $D$ . It learns to reverse the diffusion process (*i.e.*, denoise) by optimizing  $\mathcal{L}_{\text{diff}}$  in Eq. 2. (b) **Learning a Policy with the Learned Diffusion Model:** we propose a diffusion model objective  $\mathcal{L}_{\text{DM}}$  for policy learning and jointly optimize it with the BC objective  $\mathcal{L}_{\text{BC}}$ . Specifically,  $\mathcal{L}_{\text{DM}}$  is computed based on processing a sampled state-action pair  $(s, a)$  and a state-action pair  $(s, \hat{a})$  with the action  $\hat{a}$  predicted by the policy  $\pi$  with  $\mathcal{L}_{\text{diff}}$ .

182 propose to view  $\mathcal{L}_{\text{diff}}(s, a, \phi)$  as an estimate of how well the state-action pair  $(s, a)$  fits the state-action  
 183 distribution that  $\phi$  learns from.

#### 184 4.2.2 Learning a Policy with Diffusion Model Loss

185 A diffusion model  $\phi$  trained on the expert distribution can produce an estimate  $\mathcal{L}_{\text{diff}}(s, a, \phi)$  indicating  
 186 how well a state-action pair  $(s, a)$  fits the expert distribution. We propose to leverage this signal to  
 187 guide a policy to imitate the expert. Specifically, given a state-action  $(s, a)$  sampled from  $D$ , the  $\pi$   
 188 predicts an action given the state  $\hat{a} \sim \pi(s)$  by optimizing

$$\mathcal{L}_{\text{diff}}^{\text{agent}} = \mathcal{L}_{\text{diff}}(s, \hat{a}, \phi) = \|\hat{\epsilon}(s, \hat{a}, n) - \epsilon\|^2. \quad (3)$$

189 Intuitively, the policy learns to predict actions that are indistinguishable from the expert actions for  
 190 the diffusion model conditioning on the same set of states.

191 We hypothesize that learning a policy to optimize Eq. 3 can be unstable, especially for state-action  
 192 pairs that are not well-modeled by the diffusion model, which yield a high value of  $\mathcal{L}_{\text{diff}}$  even with  
 193 expert state-action pairs. Therefore, we propose to normalize the agent diffusion loss  $\mathcal{L}_{\text{diff}}^{\text{agent}}$  with an  
 194 expert diffusion loss  $\mathcal{L}_{\text{diff}}^{\text{expert}}$ , which can be computed with expert state-action pairs  $(s, a)$  as follows:

$$\mathcal{L}_{\text{diff}}^{\text{expert}} = \mathcal{L}_{\text{diff}}(s, a, \phi) = \|\hat{\epsilon}(s, a, n) - \epsilon\|^2. \quad (4)$$

195 We propose to optimize the diffusion model loss  $\mathcal{L}_{\text{DM}}$  based on calculating the difference between  
 196 the above agent and expert diffusion losses:

$$\mathcal{L}_{\text{DM}} = \max(\mathcal{L}_{\text{diff}}^{\text{agent}} - \mathcal{L}_{\text{diff}}^{\text{expert}}, 0). \quad (5)$$

#### 197 4.3 Combining the Two Objectives

198 Our goal is to learn a policy that benefits from both modeling the conditional probability and the joint  
 199 probability of expert behaviors. To this end, we propose to augment a BC policy that optimizes the  
 200 BC loss  $L_{\text{BC}}$  in Eq. 1 by jointing optimizing the proposed diffusion model loss  $L_{\text{DM}}$  in Eq. 5, which  
 201 encourages the policy to predict actions that fit the expert joint probability captured by a diffusion  
 202 model. To learn from both the BC loss and the diffusion model loss, we train the policy to optimize

$$\mathcal{L}_{\text{total}} = \mathcal{L}_{\text{BC}} + \lambda \mathcal{L}_{\text{DM}}, \quad (6)$$

203 where  $\lambda$  is a coefficient that determines the importance of the diffusion model loss relative to the BC  
 204 loss. We analyze the effect of the coefficient in Section 5.6.1.

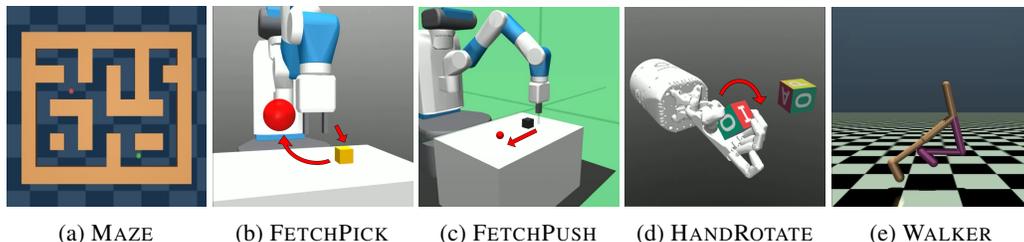


Figure 3: **Environments & Tasks.** (a) **MAZE**: A point-mass agent (green) in a 2D maze learns to navigate from its start location to a goal location (red). (b)-(c) **FETCHPICK and FETCHPUSH**: The robot arm manipulation tasks employ a 7-DoF Fetch robotics arm. **FETCHPICK** requires picking up an object (yellow cube) from the table and moving it to a target location (red); **FETCHPUSH** requires the arm to push an object (black cube) to a target location (red). (d) **HANDROTATE**: This dexterous manipulation task requires a Shadow Dexterous Hand to in-hand rotate a block to a target orientation. (e) **WALKER**: This locomotion task requires learning a bipedal walker policy to walk as fast as possible while maintaining its balance.

## 205 5 Experiments

206 We design experiments in various continuous control domains, including navigation, robot arm  
 207 manipulation, dexterous manipulation, and locomotion, to compare our proposed framework (DBC)  
 208 to its variants and baselines.

### 209 5.1 Experimental Setup

210 This section describes the environments, tasks, and expert demonstrations used for learning and  
 211 evaluation. More details can be found in Section A.

212 **Navigation.** To evaluate our method on a navigation task, we choose MAZE, a maze environment  
 213 proposed in Fu et al. [2020] (maze2d-medium-v2), as illustrated in Figure 3a. This task features  
 214 a point-mass agent in a 2D maze learning to navigate from its start location to a goal location by  
 215 iteratively predicting its  $x$  and  $y$  acceleration. The agent’s beginning and final locations are chosen  
 216 randomly. We collect 100 demonstrations with 18,525 transitions using a controller.

217 **Robot Arm Manipulation.** We evaluate our method in a robot arm manipulation domain with  
 218 two 7-DoF Fetch tasks: **FETCHPICK** and **FETCHPUSH**, as illustrated in Figure 3c and Figure  
 219 3b. **FETCHPICK** requires picking up an object from the table and lifting it to a target location;  
 220 **FETCHPUSH** requires the arm to push an object to a target location. We use the demonstrations  
 221 provided in Lee et al. [2021] for these tasks. Each dataset contains 10k transitions (303 trajectories  
 222 for **FETCHPICK** and 185 trajectories for **FETCHPUSH**).

223 **Dexterous Manipulation.** In **HANDROTATE**, we further evaluate our method on a challenging  
 224 environment proposed in Plappert et al. [2018], where a 24-DoF Shadow Dexterous Hand learns  
 225 to in-hand rotate a block to a target orientation, as illustrated in Figure 3d. This environment has  
 226 a high-dimensional state space (68D) and action space (20D). We collected 10k transitions (515  
 227 trajectories) from a SAC [Haarnoja et al., 2018] expert policy trained for 10M environment steps.

228 **Locomotion.** For locomotion, we leverage the **WALKER** environment Brockman et al. [2016], which  
 229 requires a bipedal agent to walk as fast as possible while maintaining its balance, as illustrated in  
 230 Figure 3e. We use the demonstrations provided by Kostrikov [2018], which contains 5 trajectories  
 231 with 5k state-action pairs.

### 232 5.2 Baselines

233 We compare our method DBC with the following baselines.

- 234 • **BC** learns to imitate an expert by modeling the conditional probability  $p(a|s)$  of the expert  
 235 behaviors via optimizing the BC loss  $\mathcal{L}_{BC}$  in Eq. 1.

Table 1: **Experimental Result.** We report the mean and the standard deviation of success rate (MAZE, FETCHPICK, FETCHPUSH, HANDROTATE) and return (WALKER), evaluated over three random seeds. Our proposed method (DBC) outperforms the baselines on MAZE, FETCHPICK, FETCHPUSH, HANDROTATE, and performs competitively against the best performing baseline on WALKER.

Method	MAZE	FETCHPICK	FETCHPUSH	HANDROTATE	WALKER
BC	79.35% $\pm$ 5.05%	69.15% $\pm$ 5.00%	66.02% $\pm$ 6.88%	55.48% $\pm$ 3.97%	<b>7066.61</b> $\pm$ 22.79
Implicit BC	81.43% $\pm$ 4.88%	72.27% $\pm$ 6.71%	77.70% $\pm$ 4.42%	14.52% $\pm$ 3.04%	685.92 $\pm$ 150.26
Diffusion Policy	73.34% $\pm$ 5.30%	74.37% $\pm$ 3.80%	86.93% $\pm$ 3.26%	58.59% $\pm$ 2.85%	6429.87 $\pm$ 356.70
DBC	<b>86.99%</b> $\pm$ 2.84%	<b>88.71%</b> $\pm$ 6.46%	<b>89.50%</b> $\pm$ 3.99%	<b>60.34%</b> $\pm$ 4.60%	7057.42 $\pm$ 36.19

- 236 • **Implicit BC (IBC)** [Florence et al., 2022] models expert state-action pairs with an energy-based  
237 model. For inference, we implement the derivative-free optimization algorithm proposed in IBC,  
238 which samples actions iteratively to select the desired action with the minimum predicted energy.  
239 This baseline serves a representative of the methods that solely model the joint probability  $p(s, a)$   
240 of the expert behaviors.
- 241 • **Diffusion policy** refers to the methods that learn a conditional diffusion model as a policy [Chi  
242 et al., 2023, Reuss et al., 2023]. Specifically, we implement this baseline based on Pearce et al.  
243 [2023]. We include this baseline to analyze the effectiveness of using diffusion models as a policy  
244 or as a learning objective (ours).

### 245 5.3 Experimental Results

246 We report the experimental results in terms of success rate (MAZE, FETCHPICK, FETCHPUSH,  
247 HANDROTATE), and return (WALKER) in Table 1. The details of model architecture can be found  
248 in Section C. Training and evaluation details can be found in Section D. Additional analysis and  
249 experimental results can be found in Section E and Section F.

250 **Overall Task Performance.** Our proposed method DBC achieves the highest success rates, out-  
251 performing our baselines in all the goal-directed tasks (MAZE, FETCHPICK, FETCHPUSH, and  
252 HANDROTATE) and perform competitively in WALKER compared to the best-performing baseline  
253 (BC). We hypothesize the improvement in the goal-directed tasks can be mostly attributed to the  
254 better generalization ability since starting positions and the goals are randomized during evaluation  
255 and therefore requires the policy to deal with unseen situation. To verify this hypothesis, we further  
256 evaluate the baselines and our method in FETCHPICK and FETCHPUSH with different levels of  
257 randomization in Section E.

258 **Locomotion.** Unlike the goal-directed tasks, we do not observe significant improvement but competi-  
259 tive results from DBC compared to the best-performing baseline (BC). We hypothesize that this is  
260 because locomotion tasks such as WALKER, with sufficient expert demonstrations and little random-  
261 ness, do not require generalization during inference. The agent can simply follow the closed-loop  
262 progress of the expert demonstrations, resulting in both BC (7066.61) and DBC (7057.42) performing  
263 similarly to the expert with an average return of 7063.72. On the other hand, we hypothesize that  
264 Diffusion Policy performs slightly worse due to its design for modeling multimodal behaviors, which  
265 is contradictory to learning from this single-mode simulated locomotion task.

266 **Action Space Dimension.** While Implicit BC models the joint distribution and generalizes better,  
267 it requires time-consuming actions sampling and optimization during inference. Moreover, such  
268 procedure may not scale well to high-dimensional action spaces. Our Implicit BC baseline with  
269 a derivative-free optimizer struggles in HANDROTATE and WALKER environments, whose action  
270 dimensions are 20 and 6, respectively. This is consistent with Florence et al. [2022], which reports  
271 that the optimizer failed to solve tasks with an action dimension larger than 5. In contrast, our  
272 proposed DBC can handle high-dimensional action spaces.

273 **Inference Efficiency.** To evaluate the inference efficiency, we measure and report the number of  
274 evaluation episodes per second ( $\uparrow$ ) for IBC (9.92), Diffusion Policy (1.38), and DBC (**30.79**) on  
275 an NVIDIA RTX 3080 Ti GPU in MAZE. This can be attributed to the fact that DBC and BC  
276 model the conditional probability  $p(a|s)$  and can directly map states to actions during inference. In  
277 contrast, Implicit BC requires action sampling and optimization, while Diffusion Policy is required to  
278 iteratively denoise sampled noises. This verifies the efficiency of modeling the conditional probability.

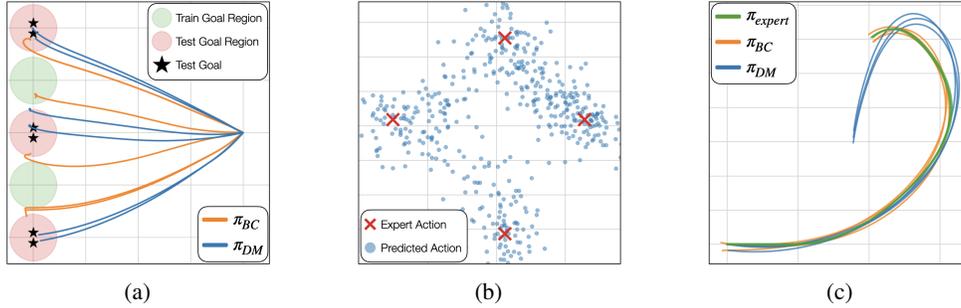


Figure 4: **Comparing Modeling Conditional Probability and Joint Probability.** (a) **Generalization.** We collect expert trajectories from a PPO policy learning to navigate to goals sampled from the green regions. Then, we learn a policy  $\pi_{BC}$  to optimize  $\mathcal{L}_{BC}$ , and another policy  $\pi_{DM}$  to optimize  $\mathcal{L}_{DM}$  with a diffusion model trained on the expert distribution. We evaluate the two policies by sampling goals from the red regions, which requires the ability to generalize.  $\pi_{BC}$  (orange) struggles at generalizing to unseen goals, whereas  $\pi_{DM}$  (blue) can generalize (*i.e.*, extrapolate) to some extent. (b)-(c) **Manifold overfitting.** We collect the green spiral trajectories from a script policy, whose actions are visualized as red crosses. We then train and evaluate  $\pi_{BC}$  and  $\pi_{DM}$ . The trajectories of  $\pi_{BC}$  (orange) can closely follow the expert trajectories (green), while the trajectories of  $\pi_{DM}$  (blue) drastically deviates from expert’s. This is because the diffusion model struggles at modeling such expert action distribution with a lower intrinsic dimension, which can be observed from poorly predicted actions (blue dots) produced by the diffusion model.

#### 279 5.4 Comparing Modeling Conditional Probability and Joint Probability

280 This section aims to empirically identify the limitations of modeling *either* the conditional *or* the  
 281 joint probability in an open maze environment implemented with [Fu et al., 2020].

282 **Generalization.** We aim to investigate if learning from the BC loss alone struggles at generalization  
 283 (*conditional*) and examine if guiding the policy using the diffusion model loss yields improved  
 284 generalization ability (*joint*). We collect trajectories of a PPO policy learning to navigate from  
 285 (5, 3) to goals sampled around (1, 2) and (1, 4) (green), as shown in Figure 4a. Given these expert  
 286 trajectories, we learn a policy  $\pi_{BC}$  to optimize Eq. 1 and another policy  $\pi_{DM}$  to optimize Eq. 5. Then,  
 287 we evaluate the two policies by sampling goals around (1, 1), (1, 3), and (1, 5) (red), which requires  
 288 the ability to generalize. Visualized trajectories of the two policies in Figure 4a show that  $\pi_{BC}$   
 289 (orange) fails to generalize to unseen goals, whereas  $\pi_{DM}$  (blue) can generalize (*i.e.*, extrapolate) to  
 290 some extent. This verifies our motivation to augment BC with the diffusion model loss.

291 **Manifold overfitting.** We aim to examine if modeling the joint probability is difficult when observed  
 292 high-dimensional data lies on a low-dimensional manifold (*i.e.*, manifold overfitting). We collect  
 293 trajectories from a script policy that executes actions (0.5, 0), (0, 0.5), (−0.7, 0), and (0, −0.7) (red  
 294 crosses in Figure 4b), each for 40 consecutive time steps, resulting the green spiral trajectories  
 295 visualized in Figure 4c.

296 Given these expert demonstrations, we learn a policy  $\pi_{BC}$  to optimize Eq. 1, and another policy  
 297  $\pi_{DM}$  to optimize Eq. 5 with a diffusion model trained on the expert distribution. Figure 4b shows  
 298 that the diffusion model struggles at modeling such expert action distribution with a lower intrinsic  
 299 dimension. As a result, Figure 4c show that the trajectories of  $\pi_{DM}$  (blue) drastically deviates from  
 300 the expert trajectories (green) as the diffusion model cannot provide effective loss. On the other hand,  
 301 the trajectories of  $\pi_{BC}$  (orange) is able to closely follow expert’s. This verifies our motivation to  
 302 complement modeling the joint probability with modeling the conditional probability (*i.e.*, BC).

#### 303 5.5 Comparing Different Generative Models

304 Our proposed framework employs a diffusion model (DM) to model the joint probability of expert  
 305 state-action pairs and utilizes it to guide policy learning. To justify our choice, we explore using other  
 306 popular generative models to replace the diffusion model in MAZE. We consider energy-based models  
 307 (EBMs) [Du and Mordatch, 2019, Song and Kingma, 2021], variational autoencoder (VAEs) [Kingma  
 308 and Welling, 2014], and generative adversarial networks (GANs) Goodfellow et al. [2014]. Each

Table 2: **Generative Models.** We compare using different generative models to model the expert distribution and guide policy learning in MAZE.

Method	without BC	with BC
BC	N/A	79.35% $\pm$ 5.05%
EBM	49.09% $\pm$ 15.15%	80.00% $\pm$ 4.06%
VAE	48.47% $\pm$ 7.57%	82.31% $\pm$ 5.84%
GAN	50.29% $\pm$ 8.27%	71.64% $\pm$ 5.50%
DM	<b>53.51%</b> $\pm$ 4.20%	<b>86.99%</b> $\pm$ 2.84%

Table 3: **Effect of  $\lambda$ .** We experiment with different values of  $\lambda$  in MAZE, each evaluated over three random seeds.

$\lambda$	Success Rate
1	85.40% $\pm$ 4.37%
2	85.64% $\pm$ 3.69%
5	<b>86.99%</b> $\pm$ 2.84%
10	85.46% $\pm$ 4.47%
20	85.17% $\pm$ 2.61%

309 generative model learns to model expert state-action pairs. To guide policy learning, given a predicted  
 310 state-action pair  $(s, \hat{a})$  we use the estimated energy of an EBM, the reconstruction error of a VAE,  
 311 and the discriminator output of a GAN to optimize a policy with or without the BC loss. Training  
 312 details can be found in Section D.3.

313 Table 2 compares using different generative models to model the expert distribution and guide  
 314 policy learning. All the generative model-guide policies can be improved by adding the BC loss,  
 315 justifying our motivation to complement modeling the joint probability with modeling the conditional  
 316 probability. With or without the BC loss, the diffusion model-guided policy achieves the best  
 317 performance compared to other generative models, verifying our choice of the generative model.

## 318 5.6 Ablation Study

319 In this section, we investigate the effect of the diffusion model loss coefficient  $\lambda$  (Section 5.6.1) and  
 320 examine the effect of the normalization term  $\mathcal{L}_{\text{diff}}^{\text{expert}}$  in the diffusion model loss  $\mathcal{L}_{\text{DM}}$  (Section 5.6.2).

### 321 5.6.1 Effect of the Diffusion Model Loss Coefficient $\lambda$

322 We examine the impact of varying the coefficient of the diffusion model loss  $\lambda$  in Eq. 6 in MAZE.  
 323 The result presented in Table 3 shows that  $\lambda = 5$  yields the best performance. A higher or lower  $\lambda$   
 324 leads to worse performance, demonstrating how modeling the conditional probability ( $\mathcal{L}_{\text{BC}}$ ) and the  
 325 joint probability ( $\mathcal{L}_{\text{DM}}$ ) can complement each other.

### 326 5.6.2 Effect of the Normalization Term $\mathcal{L}_{\text{diff}}^{\text{expert}}$

327 We aim to investigate whether normalizing the diffusion model loss  $\mathcal{L}_{\text{DM}}$  with the expert diffusion  
 328 model loss  $\mathcal{L}_{\text{diff}}^{\text{expert}}$  yields improved performance in MAZE. We train a variant of DBC where only  
 329  $\mathcal{L}_{\text{diff}}^{\text{agent}}$  in Eq. 3 instead of  $\mathcal{L}_{\text{DM}}$  in Eq. 5 is used to augment BC. This variant learning from an  
 330 unnormalized diffusion model loss achieves an average success rate of 80.20%, worse than the full  
 331 DBC (86.99%). This justifies the effectiveness of the proposed normalization term  $\mathcal{L}_{\text{diff}}^{\text{expert}}$  in  $\mathcal{L}_{\text{DM}}$ .

## 332 6 Conclusion

333 We propose an imitation learning framework that benefits from modeling both the conditional  
 334 probability  $p(a|s)$  and the joint probability  $p(s, a)$  of the expert distribution. Our proposed diffusion  
 335 model-augmented behavioral cloning (DBC) employs a diffusion model trained to model expert  
 336 behaviors and learns a policy to optimize both the BC loss and our proposed diffusion model loss.  
 337 Specifically, the BC loss captures the conditional probability  $p(a|s)$  from expert state-action pairs,  
 338 which directly guides the policy to replicate the expert’s action. On the other hand, the diffusion  
 339 model loss models the joint distribution of expert’s state-action pairs  $p(s, a)$ , which provides an  
 340 evaluation of how well the predicted action aligned with the expert distribution. DBC outperforms  
 341 baselines or achieves competitive performance in various continuous control tasks in navigation,  
 342 robot arm manipulation, dexterous manipulation, and locomotion. We design additional experiments  
 343 to verify the limitations of modeling either the conditional probability or the joint probability of the  
 344 expert distribution as well as compare different generative models. Ablation studies investigate the  
 345 effect of hyperparameters and justify the effectiveness of our design choices.

346 **References**

- 347 Volodymyr Mnih, Koray Kavukcuoglu, David Silver, Andrei A Rusu, Joel Veness, Marc G Bellemare,  
348 Alex Graves, Martin Riedmiller, Andreas K Fidjeland, Georg Ostrovski, et al. Human-level control  
349 through deep reinforcement learning. *Nature*, 2015.
- 350 Timothy P. Lillicrap, Jonathan J. Hunt, Alexander Pritzel, Nicolas Heess, Tom Erez, Yuval Tassa,  
351 David Silver, and Daan Wierstra. Continuous control with deep reinforcement learning. In  
352 *International Conference on Learning Representations*, 2016.
- 353 Kai Arulkumaran, Marc Peter Deisenroth, Miles Brundage, and Anil Anthony Bharath. Deep  
354 reinforcement learning: A brief survey. *IEEE Signal Processing Magazine*, 2017.
- 355 Paul F Christiano, Jan Leike, Tom Brown, Miljan Martic, Shane Legg, and Dario Amodei. Deep  
356 reinforcement learning from human preferences. In *Advances in Neural Information Processing*  
357 *Systems*, 2017.
- 358 Jan Leike, David Krueger, Tom Everitt, Miljan Martic, Vishal Maini, and Shane Legg. Scalable agent  
359 alignment via reward modeling: a research direction. *arXiv preprint arXiv:1811.07871*, 2018.
- 360 Youngwoon Lee, Shao-Hua Sun, Sriram Somasundaram, Edward S. Hu, and Joseph J. Lim. Compos-  
361 ing complex skills by learning transition policies. In *Proceedings of International Conference on*  
362 *Learning Representations*, 2019.
- 363 Javier Garcia and Fernando Fernández. A comprehensive survey on safe reinforcement learning.  
364 *Journal of Machine Learning Research*, 2015.
- 365 Sergey Levine, Aviral Kumar, George Tucker, and Justin Fu. Offline reinforcement learning: Tutorial,  
366 review, and perspectives on open problems. *arXiv preprint arXiv:2005.01643*, 2020.
- 367 Stefan Schaal. Learning from demonstration. In *Advances in Neural Information Processing Systems*,  
368 1997.
- 369 Takayuki Osa, Joni Pajarinen, Gerhard Neumann, J Andrew Bagnell, Pieter Abbeel, Jan Peters, et al.  
370 An algorithmic perspective on imitation learning. *Foundations and Trends® in Robotics*, 2018.
- 371 Jonathan Ho and Stefano Ermon. Generative adversarial imitation learning. In *Advances in Neural*  
372 *Information Processing Systems*, 2016.
- 373 Konrad Zolna, Scott Reed, Alexander Novikov, Sergio Gomez Colmenarejo, David Budden, Serkan  
374 Cabi, Misha Denil, Nando de Freitas, and Ziyu Wang. Task-relevant adversarial imitation learning.  
375 In *Conference on Robot Learning*, 2021.
- 376 Ilya Kostrikov, Kumar Krishna Agrawal, Debidatta Dwibedi, Sergey Levine, and Jonathan Tompson.  
377 Discriminator-actor-critic: Addressing sample inefficiency and reward bias in adversarial imitation  
378 learning. In *International Conference on Learning Representations*, 2019.
- 379 Andrew Y. Ng and Stuart J. Russell. Algorithms for inverse reinforcement learning. In *International*  
380 *Conference on Machine Learning*, 2000.
- 381 Pieter Abbeel and Andrew Y Ng. Apprenticeship learning via inverse reinforcement learning. In  
382 *International Conference on Machine Learning*, 2004.
- 383 Dean A Pomerleau. Alvin: An autonomous land vehicle in a neural network. In *Advances in Neural*  
384 *Information Processing Systems*, 1989.
- 385 Michael Bain and Claude Sammut. A framework for behavioural cloning. In *Machine Intelligence*  
386 *15*, 1995.
- 387 Tung Nguyen, Qinqing Zheng, and Aditya Grover. Reliable conditioning of behavioral cloning for  
388 offline reinforcement learning. *arXiv preprint arXiv:2210.05158*, 2023.
- 389 Pete Florence, Corey Lynch, Andy Zeng, Oscar A Ramirez, Ayzaan Wahid, Laura Downs, Adrian  
390 Wong, Johnny Lee, Igor Mordatch, and Jonathan Tompson. Implicit behavioral cloning. In  
391 *Conference on Robot Learning*, 2022.

- 392 Aditya Ganapathi, Pete Florence, Jake Varley, Kaylee Burns, Ken Goldberg, and Andy Zeng. Implicit  
393 kinematic policies: Unifying joint and cartesian action spaces in end-to-end robot learning. In  
394 *International Conference on Robotics and Automation*, 2022.
- 395 Abdoulaye O Ly and Moulay Akhloufi. Learning to drive by imitation: An overview of deep behavior  
396 cloning methods. *IEEE Transactions on Intelligent Vehicles*, 2020.
- 397 Jack Harmer, Linus Gisslén, Jorge del Val, Henrik Holst, Joakim Bergdahl, Tom Olsson, Kristoffer  
398 Sjöo, and Magnus Nordin. Imitation learning with concurrent actions in 3d games. In *IEEE*  
399 *Conference on Computational Intelligence and Games*, 2018.
- 400 Faraz Torabi, Garrett Warnell, and Peter Stone. Behavioral cloning from observation. In *International*  
401 *Joint Conference on Artificial Intelligence*, 2018.
- 402 Stéphane Ross, Geoffrey Gordon, and Drew Bagnell. A reduction of imitation learning and structured  
403 prediction to no-regret online learning. In *International Conference on Artificial Intelligence and*  
404 *Statistics*, 2011.
- 405 Faraz Torabi, Garrett Warnell, and Peter Stone. Generative adversarial imitation from observation.  
406 *ICML*, 2019.
- 407 Ian Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair,  
408 Aaron Courville, and Yoshua Bengio. Generative adversarial nets. In *Advances in Neural*  
409 *Information Processing Systems*, 2014.
- 410 Minshuo Chen, Yizhou Wang, Tianyi Liu, Zhuoran Yang, Xingguo Li, Zhaoran Wang, and Tuo Zhao.  
411 On computation and generalization of generative adversarial imitation learning. In *International*  
412 *Conference on Learning Representations*, 2020.
- 413 Justin Fu, Katie Luo, and Sergey Levine. Learning robust rewards with adversarial inverse reinforce-  
414 ment learning. In *International Conference on Learning Representations*, 2018.
- 415 Youngwoon Lee, Andrew Szot, Shao-Hua Sun, and Joseph J. Lim. Generalizable imitation learning  
416 from observation via inferring goal proximity. In *Neural Information Processing Systems*, 2021.
- 417 Tim Pearce, Tabish Rashid, Anssi Kanervisto, David Bignell, Mingfei Sun, Raluca Georgescu,  
418 Sergio Valcarcel Macua, Shan Zheng Tan, Ida Momennejad, Katja Hofmann, and Sam Devlin.  
419 Imitating human behaviour with diffusion models. In *International Conference on Learning*  
420 *Representations*, 2023.
- 421 Cheng Chi, Siyuan Feng, Yilun Du, Zhenjia Xu, Eric Cousineau, Benjamin Burchfiel, and Shuran  
422 Song. Diffusion policy: Visuomotor policy learning via action diffusion. In *Robotics: Science and*  
423 *Systems*, 2023.
- 424 Moritz Reuss, Maximilian Li, Xiaogang Jia, and Rudolf Lioutikov. Goal-conditioned imitation  
425 learning using score-based diffusion policies. *arXiv preprint arXiv:2304.02532*, 2023.
- 426 Felipe Codevilla, Eder Santana, Antonio M López, and Adrien Gaidon. Exploring the limitations of  
427 behavior cloning for autonomous driving. In *International Conference on Computer Vision*, 2019.
- 428 Lingguang Wang, Carlos Fernandez, and Christoph Stiller. High-level decision making for automated  
429 highway driving via behavior cloning. *IEEE Transactions on Intelligent Vehicles*, 2022.
- 430 Yilun Du and Igor Mordatch. Implicit generation and modeling with energy based models. In *Neural*  
431 *Information Processing Systems*, 2019.
- 432 Yang Song and Diederik P Kingma. How to train your energy-based models. *arXiv preprint*  
433 *arXiv:2101.03288*, 2021.
- 434 Diederik P Kingma and Max Welling. Auto-encoding variational bayes. In *International Conference*  
435 *on Learning Representations*, 2014.
- 436 Danilo Rezende and Shakir Mohamed. Variational inference with normalizing flows. In *International*  
437 *Conference on Machine Learning*, 2015.

- 438 Laurent Dinh, Jascha Sohl-Dickstein, and Samy Bengio. Density estimation using real NVP. In  
439 *International Conference on Learning Representations*, 2017.
- 440 Qitian Wu, Rui Gao, and Hongyuan Zha. Bridging explicit and implicit deep generative models via  
441 neural stein estimators. In *Neural Information Processing Systems*, 2021.
- 442 Gabriel Loaiza-Ganem, Brendan Leigh Ross, Jesse C Cresswell, and Anthony L Caterini. Diagnosing  
443 and fixing manifold overfitting in deep generative models. *Transactions on Machine Learning*  
444 *Research*, 2022.
- 445 Jascha Sohl-Dickstein, Eric Weiss, Niru Maheswaranathan, and Surya Ganguli. Deep unsupervised  
446 learning using nonequilibrium thermodynamics. In *International Conference on Machine Learning*.  
447 PMLR, 2015.
- 448 Alexander Quinn Nichol and Prafulla Dhariwal. Improved denoising diffusion probabilistic models.  
449 In *International Conference on Machine Learning*, 2021.
- 450 Prafulla Dhariwal and Alexander Nichol. Diffusion models beat gans on image synthesis. *Neural*  
451 *Information Processing Systems*, 2021.
- 452 A Jain J Ho. Denoising diffusion probabilistic models. In *Advances in Neural Information Processing*  
453 *Systems*, 2020.
- 454 Justin Fu, Aviral Kumar, Ofir Nachum, George Tucker, and Sergey Levine. D4rl: Datasets for deep  
455 data-driven reinforcement learning. *arXiv preprint arXiv:2004.07219*, 2020.
- 456 Matthias Plappert, Marcin Andrychowicz, Alex Ray, Bob McGrew, Bowen Baker, Glenn Powell,  
457 Jonas Schneider, Josh Tobin, Maciek Chociej, Peter Welinder, et al. Multi-goal reinforcement learn-  
458 ing: Challenging robotics environments and request for research. *arXiv preprint arXiv:1802.09464*,  
459 2018.
- 460 Tuomas Haarnoja, Aurick Zhou, Pieter Abbeel, and Sergey Levine. Soft actor-critic: Off-policy max-  
461 imum entropy deep reinforcement learning with a stochastic actor. In *Proceedings of International*  
462 *Conference on Machine Learning (ICML)*, 2018.
- 463 Greg Brockman, Vicki Cheung, Ludwig Pettersson, Jonas Schneider, John Schulman, Jie Tang, and  
464 Wojciech Zaremba. Openai gym, 2016.
- 465 Ilya Kostrikov. Pytorch implementations of reinforcement learning algorithms. <https://github.com/ikostrikov/pytorch-a2c-ppo-acktr-gail>, 2018.  
466