

1	Contents	
2	A Foundation of Diffusion Models	2
3	A.1 SDEs/ODEs and Solvers	2
4	A.2 Guided Sampling Methods	2
5	B Related Works	3
6	C Details of Experimental Setup	3
7	C.1 Offline Reinforcement Learning Environments and Datasets	3
8	C.2 Offline Imitation Learning Environments and Datasets	4
9	D Additional Experiments	5
10	D.1 Impact of Model Size in RL Benchmarks	5
11	D.2 Impact of Diffusion Backbones and Sampling Steps (Full Results)	6
12	D.3 Additional Analyses of DMs in IL Benchmarks	7
13	E Experimental Details	7
14	E.1 Computing Resources	7
15	E.2 Evaluation Metrics	8
16	E.3 Algorithm Hyperparameters	8
17	F Implemented Diffusion Models	8
18	F.1 DDPM/DDIM/DPM-Solver/DPM-Solver++	8
19	F.2 EDM	10
20	F.3 Rectified Flow	11
21	G Implemented Algorithms	11
22	G.1 Diffusion Planners	11
23	G.2 Diffusion Policies	12
24	G.3 Diffusion Data Synthesizers.	13
25	H Limitations, Challenges, and Future Directions	13
26	I Potential Social Impact	14
27	J License	14

28 A Foundation of Diffusion Models

29 A.1 SDEs/ODEs and Solvers

30 Assume a D -dimensional random variable $\mathbf{x}_0 \sim \mathbb{R}^D$ with an unknown distribution $q_0(\mathbf{x}_0)$ ¹. Diffu-
 31 sion Models (DMs) [24, 47] define a *forward process* $\{\mathbf{x}_t\}_{t \in [0, T]}$ with $T > 0$ by the *noise schedule*
 32 $\{\alpha_t, \sigma_t\}_{t \in [0, T]}$, such that $\forall t \in [0, T]$, \mathbf{x}_t satisfies

$$\mathbf{x}_t = \alpha_t \mathbf{x}_0 + \sigma_t \boldsymbol{\epsilon}, \quad \boldsymbol{\epsilon} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}), \quad (1)$$

33 where $\alpha_t, \sigma_t \in \mathbb{R}^+$ are differentiable functions of t and the *signal-to-noise-ratio* (SNR) α_t^2/σ_t^2 is
 34 strictly decreasing w.r.t t . The forward process in Equation (1) can also be described as a stochastic
 35 differential equation (SDE) for any $t \in [0, T]$ [24]:

$$d\mathbf{x}_t = f(t)\mathbf{x}_t dt + g(t)d\mathbf{w}_t, \quad \mathbf{x}_0 \sim q_0(\mathbf{x}_0), \quad (2)$$

36 where $\mathbf{w}_t \in \mathbb{R}^D$ is the standard Wiener process, and $f(t) = \frac{d \log \alpha_t}{dt}$, $g^2(t) = \frac{d\sigma_t^2}{dt} - 2\sigma_t^2 \frac{d \log \alpha_t}{dt}$. The
 37 SDE forward process in Equation (2) has an equivalent *reverse process* from time T to 0 [47]:

$$d\mathbf{x}_t = [f(t)\mathbf{x}_t - g^2(t)\nabla_{\mathbf{x}} \log q_t(\mathbf{x}_t)]dt + g(t)d\bar{\mathbf{w}}_t, \quad \mathbf{x}_T \sim q_T(\mathbf{x}_T), \quad (3)$$

38 where $\bar{\mathbf{w}}_t$ is a standard Wiener process in the reverse time. One can sample $q_0(\mathbf{x}_0)$ by directly solving
 39 the SDE in ??, in which the only unknown term is the *score function* $\nabla_{\mathbf{x}} \log q_t(\mathbf{x}_t)$. In practice, a
 40 neural network $\boldsymbol{\epsilon}_\theta(\mathbf{x}_t)$ parameterized by θ can be trained to approximate the scaled score function
 41 $-\sigma_t \nabla_{\mathbf{x}} \log q_t(\mathbf{x}_t)$ by minimizing the score matching loss [19, 46, 47]:

$$\mathcal{L}(\theta) := \mathbb{E}_{t \sim \text{Uniform}(0, T), \mathbf{x}_t \sim q_t(\mathbf{x}_t)} [\|\boldsymbol{\epsilon}_\theta(\mathbf{x}_t, t) + \sigma_t \nabla_{\mathbf{x}} \log q_t(\mathbf{x}_t)\|_2^2] \quad (4)$$

$$= \mathbb{E}_{t \sim \text{Uniform}(0, T), \mathbf{x}_0 \sim q_0(\mathbf{x}_0), \boldsymbol{\epsilon} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})} [\|\boldsymbol{\epsilon}_\theta(\mathbf{x}_t, t) - \boldsymbol{\epsilon}\|_2^2]. \quad (5)$$

42 Since $\boldsymbol{\epsilon}_\theta(\mathbf{x}_t, t)$ can be considered as a predicted Gaussian noise added to \mathbf{x}_t , it is usually called
 43 the *noise prediction model*. With a well-trained noise prediction model, SDE in ?? can be solved
 44 using numerical solvers, and DDPM [19] is one such method. However, numerical solvers require
 45 discretization from T to 0, in which the randomness of the Wiener process limits the step size [25].
 46 For faster sampling, one can solve the following *probability flow ODE*, which is proven to have the
 47 same marginal distribution as that of the SDE for any $t \in [0, T]$ [47]:

$$\frac{d\mathbf{x}_t}{dt} = f(t)\mathbf{x}_t - \frac{1}{2}g^2(t)\nabla_{\mathbf{x}} \log q_t(\mathbf{x}_t), \quad \mathbf{x}_T \sim q_T(\mathbf{x}_T). \quad (6)$$

48 DDIM [46] discretizes the ODE to the first order for solving, achieving almost no loss in quality with
 49 fewer sampling steps. DPM-Solver [34, 35] leverages the semi-linearity of diffusion ODEs in ?? for
 50 exact solutions, eliminating errors in the linear terms, resulting in a higher sample quality. Some
 51 works also reformulate the framework. EDM [23] optimizes the design choices from a perspective of
 52 noise schedule and uses a specially designed score function preconditioning to improve the sample
 53 quality. Rectified flow [33], on the other hand, designs a straight probability flow ODE from the
 54 optimal transport (OT) perspective, which can straighten itself through *reflow* procedure. The straight
 55 property of Rectified flow allows high-quality generation in very few sampling steps.

56 A.2 Guided Sampling Methods

57 Guided sampling methods aim to draw samples from $q_0(\mathbf{x}_0|y)$ to generate outputs with the charac-
 58 teristics of the label y . Depending on whether an additional classifier needs to be trained, guided
 59 sampling methods are divided into two categories: classifier guidance (CG) [5] and classifier-free
 60 guidance (CFG) [20].

61 **Classifier Guidance:** For conditional sampling, the score function needs to be changed to
 62 $\nabla_{\mathbf{x}} \log q_t(\mathbf{x}_t|y)$, which can be decomposed with the Bayes Theorem:

$$\nabla_{\mathbf{x}} \log q_t(\mathbf{x}_t|y) = \nabla_{\mathbf{x}} \log q_t(\mathbf{x}_t) + \nabla_{\mathbf{x}} \log q_t(y|\mathbf{x}_t), \quad (7)$$

¹To ensure clarity, we establish the convention that the subscript t denotes the timestep in the diffusion process, while the superscript t represents the timestep in sequential decision-making problem.

where the first term can be approximated by the noise prediction model, and the second term is a *noising classifier* that predicts the label y of the corrupt data \mathbf{x}_t . In practice, an additional neural network $\mathcal{C}_\phi(\mathbf{x}_t, t, \mathbf{y})$ is trained to approximate $\log q_t(\mathbf{y}|\mathbf{x}_t)$, and its gradient is computed to guide sampling process:

$$\bar{\epsilon}_\theta(\mathbf{x}_t, t, \mathbf{y}) = \epsilon_\theta(\mathbf{x}_t, t) - w\sigma_t \nabla_{\mathbf{x}} \mathcal{C}_\phi(\mathbf{x}_t, t, \mathbf{y}), \quad (8)$$

where w stands for the guidance scale. A larger value of w sharpens the classifier, amplifying the influence of the label y .

Classifier-free Guidance: According to Equation (7), the gradient of the classifier $\nabla_{\mathbf{x}} \log q_t(\mathbf{y}|\mathbf{x}_t)$ can be written to $\nabla_{\mathbf{x}} \log q_t(\mathbf{x}_t|\mathbf{y}) - \nabla_{\mathbf{x}} \log q_t(\mathbf{x}_t)$. By training a conditional noise prediction model $\epsilon_\theta(\mathbf{x}_t, t, \mathbf{y})$, the sampling process can be guided with no additional classifier:

$$\bar{\epsilon}_\theta(\mathbf{x}_t, t, \mathbf{y}) = \epsilon_\theta(\mathbf{x}_t, t) - w\sigma_t \nabla_{\mathbf{x}} \log q_t(\mathbf{y}|\mathbf{x}_t) = \epsilon_\theta(\mathbf{x}_t, t) + w(\epsilon_\theta(\mathbf{x}_t, t, \mathbf{y}) - \epsilon_\theta(\mathbf{x}_t, t)) \quad (9)$$

where $\epsilon_\theta(\mathbf{x}_t, t) = \epsilon_\theta(\mathbf{x}_t, t, \Phi)$ is approximated by the noise prediction model conditioned on a pre-specified label Φ standing for non-conditioning. Although CFG can generate trajectories specific to condition \mathbf{y} , it may cause the agent to reject higher likelihood trajectories in sequential environments, resulting in a performance drop [40]. Therefore, some methods [4, 40, 50] set the guidance weight w to 1, i.e., no guidance paradigm.

B Related Works

In recent years, DMs have demonstrated promising performance in various domains [51, 45, 32, 24], giving rise to several high-quality DM libraries, such as Diffusers [49] and Stable Diffusion [44]. These open-source libraries have significantly promoted research and applications in related fields. However, unfortunately, these libraries are designed for multimedia such as image, audio, and video generation, lacking adaptation for decision-making tasks. This is likely because DMs play diverse roles in decision-making, with various usage patterns and many unique mechanism incorporations, creating a gap in the multimedia generation paradigm. A library specially designed for decision-making is currently missing, and most research codebases are inherited from a few pioneering studies [21, 50, 4]. While effective, their algorithm-specific mechanisms and tightly coupled system architecture make it challenging for customized development.

CleanDiffuser aims to provide an "easy-to-hack" starter kit for research needs, offering researchers more exploration possibilities. We draw from the experience of many open-source decision-making libraries. For example, we emulate stable-baselines3 [43] to carefully reproduce results to provide practitioners with reliable baselines for method comparison. However, we inject more modular design to encourage users to freely design and modify. We also follow CORL [48] in designing clean and logically clear pipelines for readability, but, considering the complexity of DMs, abandon the one-file-from-scratch approach and opt for a one-file pipeline approach to offer rich examples of how to utilize CleanDiffuser building blocks to implement decision-making algorithms. Additionally, we follow Ray [30] in providing ample parameter selection interfaces within modules, making it easy for users unfamiliar with the internal implementation to customize effortlessly. In summary, CleanDiffuser is not only the first open-sourced modularized DM library tailored for decision-making algorithms but also a new library that draws on the advanced experiences of many open-source decision-making libraries.

C Details of Experimental Setup

C.1 Offline Reinforcement Learning Environments and Datasets

We evaluate 7 diffusion-based RL algorithms implemented with CleanDiffuser on 15 offline RL tasks from 3 benchmarks, including locomotion, manipulation, and navigation. These tasks are widely recognized and extensively used in offline RL settings [27, 10, 26, 16, 50, 22, 21, 1, 7, 29, 18], enjoying significant acceptance within the research community. Visualization of these tasks is

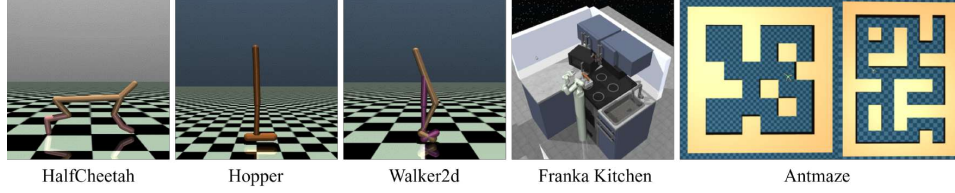


Figure 1: Visualization of Offline Reinforcement Learning Environments.

presented in Figure 1. These tasks come from the D4RL benchmark, in which the datasets are licensed under the Creative Commons Attribution 4.0 License (CC BY), and the code is licensed under the Apache 2.0 License.

Gym-MuJoCo [2] consists of three popular offline RL locomotion tasks (HalfCheetah, Hopper, Walker2d), which require controlling three Mujoco robots to achieve maximum movement speed while minimizing energy consumption under stable conditions. D4RL [9] benchmark provides three different quality levels of offline datasets: “medium” containing demonstrations of medium-level performance; “medium-replay” containing all recordings in the replay buffer observed during training until the policy reaches “medium” performance; and “medium-expert” which combines “medium” and “expert” level performance equally.

Franka Kitchen [11] requires controlling a realistic 9-DoF Franka robot arm to complete several household tasks in a kitchen environment. Algorithms are trained on “partial” and “mixed” datasets. The “partial” and “mixed” datasets consist of undirected data, where the robot performs subtasks that are not necessarily related to the goal configuration. In the “partial” dataset, a subset of the dataset is guaranteed to solve the task, meaning an imitation learning agent may learn by selectively choosing the right subsets of the data. The “mixed” dataset contains no trajectories that solve the task completely, and the RL agent must learn to assemble the relevant sub-trajectories. This dataset requires the highest degree of generalization in order to succeed.

Antmaze [9] requires controlling the 8-DoF “Ant” quadruped robot to complete maze navigation tasks. In the offline dataset, the robot only receives a reward upon reaching the goal, and the dataset contains many trajectory segments that do not lead to the endpoint, making it a difficult decision task with sparse rewards and a long horizon. The success rate of reaching the endpoint is used as the evaluation score, and common offline RL algorithms often struggle to achieve good performance.

C.2 Offline Imitation Learning Environments and Datasets

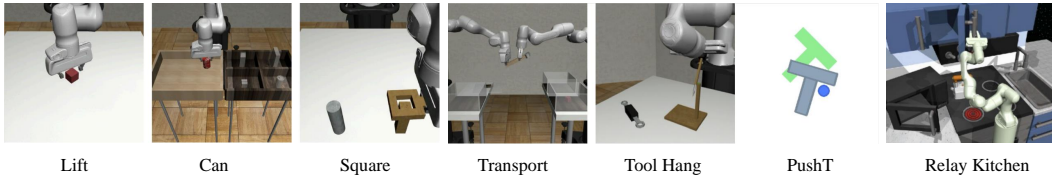


Figure 2: Visualization of Offline Imitation Learning Environments.

We evaluate 2 diffusion-based IL algorithms implemented with CleanDiffuser on 22 imitation learning tasks from 4 benchmarks, with both state and image-based observation inputs. Among them, Relay Kitchen and Robomimic support both velocity and position control. Each algorithm is trained with its best-performing action space. We provide task summary in Table 1, visualization in Figure 2, and more details below:

PushT [8] requires pushing a T-shaped block (gray) to a fixed target (red) with a circular end-effector. The task requires exploiting complex and contact-rich object dynamics to push the T block precisely, using point contacts. In this paper, we used three variants. “PushT” env has a five-dimensional state space, including the proprioception for end-effector location (agent_x , agent_y) and the

xy coordinates and angles of the blocks (block_x, block_y, block_angle). “PushT-keypoints” env includes nine 2D key points obtained from the T-block’s ground truth attitude and proprioception for end-effector location. “PushT-image” env observes the end-effector location and the top view of the RGB image. This benchmark is licensed under the Apache-2.0 License.

Relay Kitchen is proposed in Relay Policy Learning [11], commonly used to evaluate imitative learning ability. The environment consists of a 9 DoF position-controlled Franka robot interacting with a kitchen scene that includes an openable microwave, four turnable oven burners, an oven light switch, a freely movable kettle, two hinged cabinets, and a sliding cabinet door. The “relay” dataset contains 566 human demonstrations, each completing four tasks in arbitrary order. The goal is to execute as many tasks as possible, regardless of order, showcasing both short-horizon and long-horizon multimodality. This benchmark is licensed under the Apache-2.0 License.

Robomimic [37] requires controlling a robot arm to complete complex manipulation tasks from a few human demonstrations. Due to the non-Markovian nature of human demonstrations and the demonstration quality variance, learning from human datasets is significantly more challenging than learning from machine-generated datasets. Proficient-Human (PH) and Multi-Human (MH) datasets are collected by humans through remote teleoperation. The PH datasets consist of 200 demonstrations collected by a single, experienced teleoperator, while the MH datasets consist of 300 demonstrations collected by 6 teleoperators of varying proficiency, each of which provided 50 demonstrations. The benchmark consists of 5 PH tasks (Lift, Can, Square, Tool_hang, Transport) and 4 MH tasks (Lift, Can, Square, Transport). Each task has both state and image-based observation inputs. This benchmark is licensed under the MIT License.

To the best of our knowledge, the datasets and benchmarks we have used do not contain personally identifiable information or offensive content in both previous works and our works.

Table 1: **Imitation Learning Task Summary.** Obs Shape represents the low dimensional state space dimension; Image Shape represents the observation resolution of multi-view images (Camera views x W x H). PH: proficient-human demonstration, MH: multi-human demonstration, Steps: max episode steps.

Task	Low Dim Tasks	Image Tasks		Action Dim	PH Demonstration	MH Demonstration	Max Steps
	Obs Shape	Obs Shape	Image Shape				
PushT	5	N/A	N/A	2	200	N/A	300
PushT-Keypoint	20	N/A	N/A	2	200	N/A	300
PushT-Image	N/A	2	1x96x96	2	200	N/A	300
Relay Kitchen	60	N/A	N/A	9	656	N/A	280
Lift	19	9	2x84x84	7	200	300	400
Can	23	9	2x84x84	7	200	300	400
Square	23	9	2x84x84	7	200	300	500
Transport	59	18	4x84x84	7	200	300	700
Tool_hang	53	9	2x240x240	7	200	N/A	700

D Additional Experiments

D.1 Impact of Model Size in RL Benchmarks

Table 2: **Impact of Model Size in RL Benchmarks.** Performance of DD and IDQL with varying model sizes. Results correspond to the mean and standard error over 150 episode seeds.

Environment	DD			IDQL		
Model Size	4M	15M	60M	1.6M	6M	25M
HalfCheetah-m	45.3 \pm 0.3	44.5 \pm 0.1	47.1 \pm 0.1	51.5 \pm 0.1	51.5 \pm 0.1	51.7 \pm 0.1
Kitchen-m	56.5 \pm 5.8	80.5 \pm 4.1	27.7 \pm 2.1	66.5 \pm 4.1	69.2 \pm 1.0	67.5 \pm 1.8
Antmaze (mp for DD, lp for IDQL)	8.0 \pm 4.3	26.0 \pm 5.9	22.7 \pm 6.6	48.7 \pm 4.7	52.0 \pm 5.7	54.0 \pm 4.3

There is a significant disparity in network model sizes used by diffusion-based decision-making algorithms. For instance, the official implementation of DD utilizes around 60M parameters [1], while Diffuser uses 4M [21], and IDQL [16] has approximately only 1.6M parameters. These works have

168 limited discussion on the impact of model size. Therefore, we aim to explore the approximate scale
 169 of parameter sizes required for diffusion-based decision-making algorithms to function effectively.
 170 In this experiment, we test DD and IDQL at three different model sizes, starting from the default
 171 parameter size used in the main experiments and gradually increasing the parameter size by four times.
 172 The performance of the algorithms is evaluated on three tasks including locomotion, manipulation,
 173 and navigation. Results are presented in Table 2. We find that, apart from the performance of DD
 174 on Kitchen-m and Antmaze-mp, increasing the model size does not lead to significant performance
 175 gains in other cases. However, even with the performance gains brought by model size, DD can not
 176 entirely catch up with the performance of IDQL, indicating that the dominant effect on performance
 177 is still primarily driven by the algorithm rather than the model size.

178 D.2 Impact of Diffusion Backbones and Sampling Steps (Full Results)

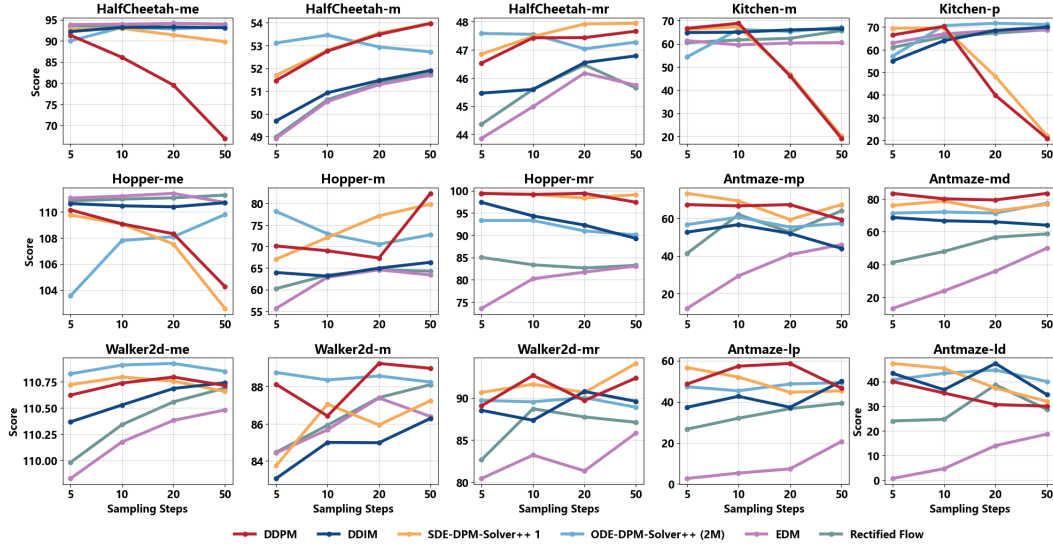


Figure 3: **Full D4RL Results of IDQL.** Performance of IDQL with various diffusion backbones and varying sampling steps. Results correspond to the mean over 150 episode seeds.

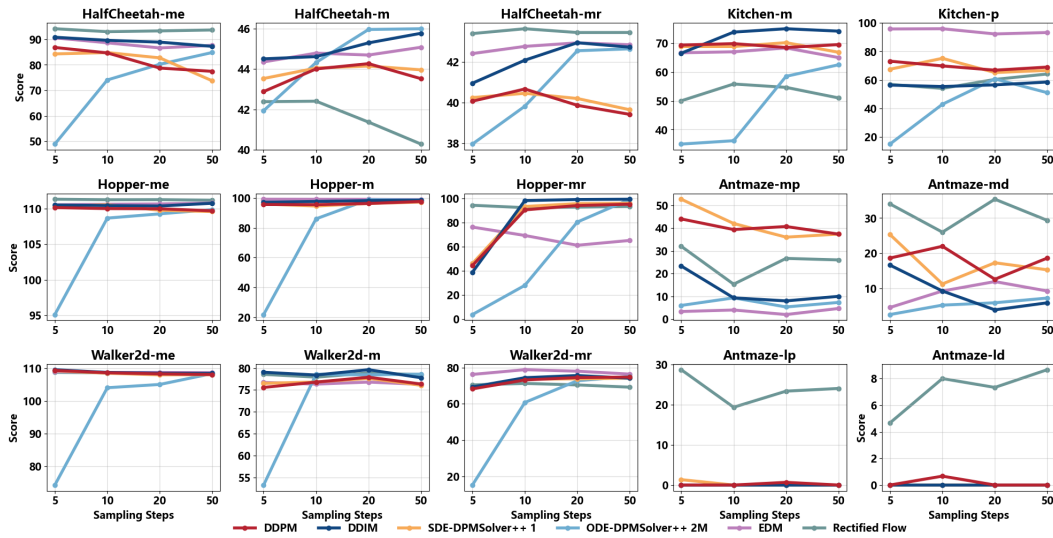


Figure 4: **Full D4RL Results of DD.** Performance of DD with various diffusion backbones and varying sampling steps. Results correspond to the mean over 150 episode seeds.

Due to space limitations in the main text, we present the full results of IDQL and DD on D4RL in Figure 3 and Figure 4. The algorithms are trained for 1×10^6 gradient steps, and the sampling steps for DD are set to 5, with other hyperparameters consistent with default settings. This experiment selects DDPM, DDIM, SDE-DPM-Solver++ 1, ODE-DPM-Solver++ (2M), EDM, and Rectified Flow as the diffusion/solver backbones. We select DDPM and DDIM because they are the first-order discretization of diffusion reverse SDE/ODE, respectively [47, 46]. We do not choose DPM-Solver because its first-order solver is equivalent to DDIM [34], and higher-order solvers may cause instability under guidance [35]. For DPM-Solver++, we select a first-order SDE solver, SDE-DPM-Solver++ 1, and a second-order ODE solver, ODE-DPM-Solver++ (2M). Since higher-order solvers can lead to instability, they are therefore not chosen. We select EDM and Rectified Flow because they have achieved excellent results in image generation but have not been widely used in the decision-making domain, to the best of our knowledge. Thanks to CleanDiffuser’s support for various solvers and varying sampling steps, the results for DDPM, DDIM, SDE-DPM-Solver++ 1, and ODE-DPM-Solver++ (2M) only require training one single model. Additionally, using different sampling steps does not require additional training. These features provide a great convenience for conducting ablation experiments. We believe these features of CleanDiffuser can also benefit future research efforts.

D.3 Additional Analyses of DMs in IL Benchmarks

Using the low-dim lift-ph task with 50 sample steps in Robomimic as a reference, we present the number of parameters and inference time for each variant of DiffusionPolicy, DiffusionBC, and ACT in table 3. Although Chi_UNet1d exhibits the best performance in many IL tasks, it has the largest model size and the slowest inference speed. Larger model size results in higher training costs, and in many real-world applications that require real-time inference, we need to make trade-offs between inference speed and performance. Compared to the transformer-based ACT algorithm, all structures of the diffusion policy exhibit slower sampling speeds because the denoising process requires multiple forwards for neural networks. This is also an important challenge that limits the application of DMs for decision-making. We also note that DiffusionBC is slower than DiffusionPolicy when using the same network architecture and model size, as DiffusionBC performs 8 additional steps of Diffusion-X sampling to mitigate OOD issues. Although the best-performing Chi_UNet1d model uses a considerable model size, simply increasing the Transformer-based DMs like DiT1d can sometimes harm performance. We discuss this in detail in appendix D.1, which is also consistent with the experimental observations of the [4]. Finding the optimal model size in applications remains an open research question.

Table 3: The Model Size and Inference Time of DiffusionPolicy and DiffusionBC in Low-Dim Lift-ph. DiffusionPolicy uses 50 sampling steps across the experiments, and DiffusionBC incorporates 8 additional Diffusion-X sampling steps.

Algorithm	Model Size (M)	Inference Time (s)
DiffusionPolicy w/ Chi_UNet1d	68.91	0.405
DiffusionPolicy w/ Chi_TFM	9.50	0.343
DiffusionPolicy w/ DiT1d	16.59	0.194
DiffusionBC w/ DiT1d	16.59	0.217
DiffusionBC w/ Pearce_MLP	0.83	0.062
ACT	7.83	0.006

E Experimental Details

E.1 Computing Resources

RL experiments are conducted on a server equipped with 2 Intel(R) Xeon(R) Gold 6326 CPUs @ 2.90GHz and 8 NVIDIA GeForce RTX3090 GPUs, and a server equipped with 2 Intel(R) Xeon(R) Gold 6326 CPUs @ 2.90GHz and 8 NVIDIA GeForce RTX2080Ti GPUs. IL experiments are conducted on a server equipped with 2 Intel(R) Xeon(R) Gold 6338 CPUs @ 2.00GHz and 8 NVIDIA A800 GPUs, and a server equipped with 2 Intel(R) Xeon(R) Gold 6338 CPUs @ 2.00GHz and 4 NVIDIA GeForce RTX3090 GPUs.

E.2 Evaluation Metrics

In the D4RL benchmark, the scores are normalized to the range between 0 and 100 with expert-normalized scores $= 100 \times \frac{\text{score} \times \text{random_score}}{\text{expert_score} - \text{random_score}}$ [9]. As for IL benchmarks, we report target area coverage as scores in the PushT benchmark and success rate in the Robomimic benchmark. In the Relay Kitchen environment, since the vast majority of human demonstrations can only complete 4 subtasks, we denote the success rate of completing the i -th subtask as p_i and report the average success rate as $\text{score} = (p_1 + p_2 + p_3 + p_4)/4$.

E.3 Algorithm Hyperparameters

Unless stated otherwise, we utilize default hyperparameters from the official implementations for most algorithms and datasets. **Configuration files and hyperparameters for each algorithm and environment are available in YAML format on our GitHub repository for reproducibility.**

Key hyperparameters for each offline RL algorithm are presented in Table 4, and each offline IL algorithm in Table 5. We also reproduce the Transformer-based ACT [52] algorithm based on the official implementation, the key hyperparameters are in Table 50000.

Table 4: Hyperparameters for Diffusion Planners, Diffusion Policies and Diffusion Data Synthesizer for RL.

Hyperparameter	Diffuser	DD	AdaptDiffuser	DQL	EDP	IDQL	SynthER
Architecture	Janner_UNet	DiT	Janner_UNet	DQL_MLP	DQL_MLP	LNResnet	LNResnet
Diffusion Model	DDPM	DDIM	DDPM	DDPM	DPM-Solver++ (2M)	DDPM	DDIM
Sampling Steps	20	20	20	5	15	5	128
Horizon	64 (Antmaze)	64 (Antmaze)	64 (Antmaze)	1	1	1	1
	32 (Otherwise)	32 (Otherwise)	32 (Otherwise)				
Temperature	0.5	0.5	0.5	0.5	0.5	0.5	1.0
Gradient Steps	1e6	1e6	1e6	2e6	2e6	2e6	1e5
Batch Size	64	64	64	256	256	256	256
Learning Rate	3e-4	3e-4	3e-4	3e-4	3e-4	3e-4	3e-4
N candidates	64	1	64	50	50	256	N/A

Table 5: Hyperparameters for DiffusionPolicy and DiffusionBC in Low-Dim and Image Tasks.

Hyperparameters	DiffusionPolicy			DiffusionBC	
Architecture	Chi_UNet1d	Chi_Transformer	DiT1d	Pearce_MLP	DiT
Diffusion Model	DDPM	DDPM	DDPM	DDPM	DDPM
Sampling Steps	5 (PushT)	5 (PushT)	5 (PushT)	50	50
	50 (Otherwise)	50 (Otherwise)	50 (Otherwise)		
Horizon	16	10	10	2	2
Obs Steps	2	2	2	2	2
Action Steps	8	8	8	1	1
Gradient Steps	1e6	1e6	1e6	1e6	1e6
Batch Size	256 (Low dim)	256 (Low dim)	256 (Low dim)	512 (Low dim)	512 (Low dim)
	64 (Image)	64 (Image)	64 (Image)	64 (Image)	64 (Image)
Temperature	1.0	1.0	1.0	1.0	1.0
Learning Rate	1e-4	1e-4	1e-4	1e-3	5e-4
Extra Sample Steps	N/A	N/A	N/A	8	8
Control Mode	Pos	Pos	Pos	Vel	Vel

F Implemented Diffusion Models

F.1 DDPM/DDIM/DPM-Solver/DPM-Solver++

Applying Solvers with One Score Function. Due to the generation processes of DDPM [19], DDIM [46], DPM-Solver [34], and DPM-Solver++ [35] can all be expressed using the same diffusion SDE/ODE [47], utilizing the same noise schedule, training just one noise predictor model enables the use of these four solvers for sampling. Recall that the diffusion ODE with noise prediction model is:

$$\frac{dx_t}{dt} = f(t)x_t + \frac{g^2(t)}{2\sigma_t} \epsilon_\theta(x_t, t). \quad (10)$$

Table 6: Hyperparameters for ACT in Low-Dim and Image Tasks.

Hyperparameters	Value
Learning Rate	1e-5
Batch Size	256 (Low dim) / 64 (Image)
# Encoder Layers	4
# Decoder Layers	7
Feedforward Dimension	256
Hidden Dimension	256
# Heads	8
Chunk size	16
Beta	10
Gradient Steps	1e6
Control Mode	Vel (Kitchen) / Pos (Otherwise)

248 Substituting $f(t) = \frac{d \log \alpha_t}{dt}$, $g^2(t) = \frac{d \sigma_t^2}{dt} - 2 \sigma_t^2 \frac{d \log \alpha_t}{dt}$, and conducting first-order discretization
 249 result in a recursive formula:

$$\mathbf{x}_t - \mathbf{x}_s = \frac{\alpha_t - \alpha_s}{\alpha_s} \mathbf{x}_s + \frac{1}{2 \sigma_s} \left[2 \sigma_s (\sigma_t - \sigma_s) - 2 \frac{\sigma_s^2}{\alpha_s} (\alpha_t - \alpha_s) \right] \epsilon_\theta(\mathbf{x}_s) \quad (11)$$

$$\mathbf{x}_t = \frac{\alpha_t}{\alpha_s} \mathbf{x}_s - \alpha_t \left(\frac{\sigma_s}{\alpha_s} - \frac{\sigma_t}{\alpha_t} \right) \epsilon_\theta(\mathbf{x}_s, s) \quad (12)$$

$$\mathbf{x}_t = \alpha_t \left(\frac{\mathbf{x}_t - \sigma_t \epsilon_\theta(\mathbf{x}_s, s)}{\alpha_s} \right) + \sqrt{\sigma_s^2} \epsilon_\theta(\mathbf{x}_s, s), \quad (13)$$

250 where t and s are the next and current sampling steps. Equation (13) is DDIM update [46]. By
 251 introduce $\beta_s = (\sigma_t / \sigma_s) \sqrt{1 - \alpha_s^2 / \alpha_t^2}$, the generative process of DDPM is:

$$\mathbf{x}_t = \alpha_t \left(\frac{\mathbf{x}_t - \sigma_t \epsilon_\theta(\mathbf{x}_s, s)}{\alpha_s} \right) + \sqrt{\sigma_s^2 - \beta_s^2} \epsilon_\theta(\mathbf{x}_s, s) + \beta_s \epsilon_s, \quad (14)$$

252 where $\epsilon_s \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ is standard Gaussian noise independent of \mathbf{x}_s . DPM-Solver leverages the
 253 semi-linearity of the diffusion ODE and formulates the exact solution by the “variation of constants”
 254 formula:

$$\mathbf{x}_t = e^{\int_s^t f(\tau) d\tau} \mathbf{x}_s + \int_s^t \left(e^{\int_\tau^t f(r) dr} \frac{g^2(\tau)}{2 \sigma_\tau} \epsilon_\theta(\mathbf{x}_\tau, \tau) \right) d\tau \quad (15)$$

$$\mathbf{x}_t = \frac{\alpha_t}{\alpha_s} \mathbf{x}_s - \alpha_t \int_s^t \frac{d\lambda_\tau}{d\tau} \frac{\sigma_\tau}{\alpha_\tau} \epsilon_\theta(\mathbf{x}_\tau, \tau) d\tau, \quad (16)$$

255 where $\lambda_t := \log(\alpha_t / \sigma_t)$ is the log-signal-to-noise-ratio (log-SNR). This formulation eliminates the
 256 approximation error of the linear term since it is exactly computed, and the non-linear term can be
 257 approximated using its Talor expansion:

$$\mathbf{x}_t = \frac{\alpha_t}{\alpha_s} \mathbf{x}_s - \alpha_t \sum_{n=0}^{k-1} \epsilon_\theta^{(n)}(\mathbf{x}, s) \int_{\lambda_s}^{\lambda_t} e^{-\lambda} \frac{(\lambda - \lambda_s)^n}{n!} d\lambda + \mathcal{O}((\lambda_t - \lambda_s)^{k+1}). \quad (17)$$

258 In CleanDiffuser, we have implemented only *DPM-Solver-1*, corresponding to the $k=1$ scenario in
 259 Equation (17), as guided sampling tends to make high-order solvers unstable [35], leading to poor
 260 performance in decision-making tasks. DPM-Solver++ alleviates this instability issue by using a
 261 data prediction model $\mathbf{x}_\theta(\mathbf{x}_t, t)$ instead of the noise prediction model $\epsilon_\theta(\mathbf{x}_t, t)$, transforming the
 262 generative process into:

$$\mathbf{x}_t = \frac{\sigma_t}{\sigma_s} \mathbf{x}_s + \sigma_t \sum_{n=0}^{k-1} \mathbf{x}_\theta^{(n)}(\mathbf{x}, s) \int_{\lambda_s}^{\lambda_t} e^\lambda \frac{(\lambda - \lambda_s)^n}{n!} d\lambda + \mathcal{O}((\lambda_t - \lambda_s)^{k+1}), \quad (18)$$

263 where $\mathbf{x}_\theta(\mathbf{x}_t, t)$ is trained to predict the original data \mathbf{x}_0 from the perturbed data \mathbf{x}_s . In
 264 CleanDiffuser, we have implemented DPM-Solver++ for $k \leq 2$, as it already yields satisfac-
 265 tory results at $k = 2$, while higher-order solvers may still lead to instability.

Although the data prediction model can mitigate the instability issue caused by guided sampling and easily clip data to address the “train-test mismatch” problem [35], there is still no definitive evidence in practice to determine the superiority of either the data prediction model or the noise prediction model. In CleanDiffuser, we provide users with the option to choose between these two prediction models and use the approximation $\mathbf{x}_t \approx \alpha_t \mathbf{x}_\theta(\mathbf{x}_t, t) + \sigma \epsilon_\theta(\mathbf{x}_t, t)$ to seamlessly switch between the two formulations to cater to the requirements of different solvers.

Noise Schedules. CleanDiffuser provides two popular noise schedules by default: *Linear Noise Schedule* [19] and *Cosine Noise Schedule* [38]. The former defines:

$$\alpha_t = \exp\left(-\frac{(\beta_1 - \beta_0)}{4}t^2 - \frac{\beta_0}{2}t\right), \quad (19)$$

where $\beta_0 = 0.1$, $\beta_1 = 20$ and $\sigma_t = \sqrt{1 - \alpha_t^2}$. The diffusion SDE/ODE is solved between $[\epsilon, T]$, where $\epsilon = 0.001$ and $T = 1$ for numerical stability. The later schedule defines:

$$\alpha_t = \frac{\cos\left(\frac{\pi}{2} \cdot \frac{t+s}{1+s}\right)}{\cos\left(\frac{\pi}{2} \cdot \frac{s}{1+s}\right)} \quad (20)$$

where $s = 0.008$ and $\sigma_t = \sqrt{1 - \alpha_t^2}$. The diffusion SDE/ODE is solved between $[\epsilon, T]$, where $\epsilon = 0.001$ and $T = 0.9946$ for numerical stability. Beyond the two schedules, CleanDiffuser allows users to fully customize new noise schedules according to the specified format to explore algorithm performance.

F.2 EDM

EDM [23] rewrites the diffusion forward process in Equation (1) as:

$$\mathbf{x}_t = s_t(\mathbf{x}_0 + \sigma_t \epsilon_t), \quad (21)$$

which can be interpreted as adding noise to a scaled version of the original data. By setting the scale $s_t \equiv 1$ to a constant, EDM obtains the following reverse process:

$$\frac{d\mathbf{x}_t}{dt} = -\dot{\sigma}_t \sigma_t \nabla_{\mathbf{x}} \log p(\mathbf{x}; \sigma_t) dt, \quad (22)$$

where $p(\mathbf{x}; \sigma_t) = p_t(\mathbf{x})$. A data prediction model $D_\theta(\mathbf{x}; \sigma)$ is trained to approximate $\mathbf{x} + \sigma^2 \nabla_{\mathbf{x}} \log p(\mathbf{x}; \sigma)$ and results in a practical generative process:

$$\mathbf{x}_t = \mathbf{x}_s + (t - s) \cdot \left(\frac{\dot{\sigma}_s}{\sigma_s} \mathbf{x}_s - \frac{\dot{\sigma}_s}{\sigma_s} D_\theta(\mathbf{x}_s; \sigma_s) \right). \quad (23)$$

One feature of EDM is that it applies preconditioning to D_θ :

$$D_\theta(\mathbf{x}; \sigma) = c_{\text{skip}}(\sigma) \mathbf{x} + c_{\text{out}}(\sigma) F_\theta(c_{\text{in}}(\sigma) \mathbf{x}; c_{\text{noise}}(\sigma)). \quad (24)$$

where F_θ is the neural network to be trained, c_{skip} modulates the skip connection, c_{in} and c_{out} scale the input and output magnitudes, and c_{noise} maps noise level σ into a conditioning input for F_θ . F_θ is trained by minimizing the noising score matching loss:

$$\mathcal{L}(\theta; \sigma) = \mathbb{E}_{\mathbf{y} \sim p_{\text{data}}, \mathbf{n} \sim \mathcal{N}(\mathbf{0}, \sigma^2 \mathbf{I})} [\lambda(\sigma) \|D_\theta(\mathbf{y} + \mathbf{n}; \sigma) - \mathbf{y}\|_2^2], \quad (25)$$

where $\lambda(\sigma)$ is the loss weight. These coefficients are optimized to achieve the following objectives: (1) inputs of F_θ have unit variance, (2) training target of F_θ have unit variance, (3) c_{skip} can minimize c_{out} so that the errors of F_θ are amplified as little as possible, and (4) the loss of F_θ has a uniform weight across noise levels. The optimization results give the following design choices: $c_{\text{skip}} = \sigma_{\text{data}}^2 / (\sigma^2 + \sigma_{\text{data}}^2)$, $c_{\text{out}} = \sigma \cdot \sigma_{\text{data}} / \sqrt{\sigma_{\text{data}}^2 + \sigma^2}$, $c_{\text{in}} = 1 / \sqrt{\sigma_{\text{data}}^2 + \sigma^2}$, $c_{\text{noise}} = \log(\sigma) / 4$, and $\lambda(\sigma) = (\sigma_{\text{data}}^2 + \sigma^2) / (\sigma_{\text{data}} \cdot \sigma)^2$.

Noise Schedule. CleanDiffuser provides only one default noise schedule, which is specially designed for EDM:

$$\sigma_t = t, \quad t \in [\sigma_{\min}, \sigma_{\max}], \quad (26)$$

where $\sigma_{\min} = 0.002$ and $\sigma_{\max} = 80$.

299 F.3 Rectified Flow

300 Rectified flow [33] is an ODE on time $t \in [0, 1]$:

$$\frac{d\mathbf{x}_t}{dt} = \mathbf{v}_\theta(\mathbf{x}_t, t), \quad (27)$$

301 where the drift force \mathbf{v}_θ is trained to drive the flow to follow the direction $(\mathbf{x}_0 - \mathbf{x}_1)$ of the linear path
302 pointing from \mathbf{x}_1 to \mathbf{x}_0 as much as possible, by solving a simple least squares regression problem:

$$\mathcal{L}(\theta) = \mathbb{E}_{\mathbf{x}_0 \sim p_0, \mathbf{x}_1 \sim p_1, t \sim \text{Uniform}(0,1)} \left[\left\| (\mathbf{x}_0 - \mathbf{x}_1) - \mathbf{v}_\theta(\mathbf{x}_t, t) \right\|_2^2 \right], \quad (28)$$

303 where $\mathbf{x}_t = t\mathbf{x}_1 + (1-t)\mathbf{x}_0$. It achieves the mutual transformation of samples from two distributions
304 p_0 and p_1 , by solving Equation (26) forward or backward. Rectified flow possesses many favorable
305 properties that allow it to continuously learn from its own sampled data to straighten the ODE flow,
306 and this procedure is called *reflow*. The straighter the ODE flow, the fewer sampling steps are needed
307 to achieve good generation quality. In an ideal scenario, if the flow becomes completely straight, then
308 we have:

$$\mathbf{x}_t = t\mathbf{x}_1 + (1-t)\mathbf{x}_0 = \mathbf{x}_1 + (1-t)\mathbf{v}(\mathbf{x}_1, 1), \quad \forall t \in [0, 1], \quad (29)$$

309 which enables one-step sampling. The Rectified Flow implemented in `CleanDiffuser` has full
310 functionality to transform samples from any two arbitrary probability distributions. By default, it
311 follows the settings in diffusion models, where p_0 is the dataset distribution and p_1 is the standard
312 Gaussian distribution.

313 G Implemented Algorithms

314 G.1 Diffusion Planners

315 **Diffuser.** [21] Diffuser is the first diffusion planning algorithm, and its paradigm has been widely
316 adopted in subsequent diffusion planning algorithms. Diffuser generates state-action pair trajectories
317 $\mathbf{x} = [x^\tau, \dots, x^{\tau+H-1}]$ from:

$$p(\mathbf{x} | \mathcal{O}^{\tau:T}) \propto p(\mathbf{x}) p(\mathcal{O}^{\tau:T} | \mathbf{x}) = p(\mathbf{x}) \prod_{t=\tau}^T \exp(r(s^t, a^t)), \quad (30)$$

318 where $\mathcal{O}^{t_1:t_2}$ is a binary random variable denoting the optimality of a trajectory from t_1 to t_2 , and \mathcal{T}
319 is the episode terminal time step of the trajectory². Therefore, it is natural to define the classifier in
320 CG as a reward function on perturbed trajectories:

$$\nabla_{\mathbf{x}} \log p_t(\mathbf{x}_t | \mathcal{O}^{\tau:T}) = \nabla_{\mathbf{x}} \log p_t(\mathbf{x}_t) + \sum_{k=\tau}^T \nabla_{s_t^k, a_t^k} r(s_t^k, a_t^k) = \nabla_{\mathbf{x}} \log p_t(\mathbf{x}_t) + \nabla_{\mathbf{x}} \mathcal{J}_\phi(\mathbf{x}_t, t), \quad (31)$$

321 where $\mathcal{J}_\phi(\mathbf{x}_t, t)$ is a neural network trained to predict the episodic cumulative reward $\sum_{k=\tau}^T r(s_t^k, a_t^k)$
322 of the perturbed trajectory \mathbf{x}_t . At each inference step, given the current state s^k , Diffuser sets
323 and freezes the first state of the trajectory as s^k and performs guided sampling in an inpainting
324 manner to generate a set of trajectories $\{\mathbf{x}_0\}$. Subsequently, it identifies the optimal trajectory
325 $\mathbf{x}_0^* = \arg \max_{\mathbf{x}_0} \mathcal{J}_\phi(\mathbf{x}_0, 0)$ that maximizes the episodic cumulative reward, and extracts the first
326 action a^k in \mathbf{x}_0^* to execute.

327 **Decision Diffuser.** [1] Decision Diffuser (DD) introduces another prominent framework that utilizes a
328 state-only trajectory formulation and implements CFG by discarding the optimality variable $\mathcal{O}^{\tau:T}$ in
329 favor of directly employing normalized episodic cumulative reward $y = \sum_{t=\tau}^T r(\mathbf{x})$ as the condition.
330 As no additional reward predictor can be used for trajectory selection, DD generates only a single

²In previous works, authors typically consider only the trajectory cumulative reward as the generative condition, i.e. using $\mathcal{O}^{\tau:\tau+H-1}$, which overlooks future optimality. Their code implementations actually use the episodic cumulative reward, i.e. $\mathcal{O}^{\tau:T}$. Therefore, we adopt this episodic cumulative reward expression.

trajectory at each inference step and employs an trained inverse dynamic model \mathcal{I}_ϕ to predict the action to be executed $a^t = \mathcal{I}_\phi(s^t, s^{t+1})$.

AdaptDiffuser. [31] Observing that the insufficient diversity of offline RL training data may limit the sample quality of DMs, AdaptDiffuser, an extension of Diffuser, proposes to utilize self-generated diverse synthetic expert data to fine-tune itself. The pipeline of AdaptDiffuser involves initially training a Diffuser as usual, then generating a large amount of synthetic expert data and using a discriminator to filter out high-quality data. Finally, fine-tuning is done on this dataset. This self-evolving process can be repeated multiple times to optimize the model, and different directions of model self-evolution can be controlled by designing different discriminators. The inference method of AdaptDiffuser is consistent with Diffuser, and its performance for seen tasks has been enhanced while also being able to adapt to unseen tasks.

G.2 Diffusion Policies

Diffusion Q-Learning. [50] Diffusion Q-learning (DQL) leverages the capability of DMs to model complex distributions, directly applying DDPM as the policy $\pi_\theta(a_0|s)$ in the RL actor-critic framework. Sampling from the policy is therefore equivalent to the denoising process of the diffusion model. The Bellman operator can be used to train the Q-value function of the diffusion policy:

$$\mathcal{L}(\phi) = \mathbb{E}_{(s^k, a^k, r, s^{k+1}) \sim \mathcal{D}, a_0^{k+1} \sim \pi_{\theta'}} \left[\left\| (r + \gamma \min_{i=1,2} Q_{\phi'_i}(s^{k+1}, a_0^{k+1})) - Q_{\phi_i}(s^k, a^k) \right\|_2^2 \right], \quad (32)$$

where ϕ_1 and ϕ_2 represent the parameters of the double Q-learning trick, ϕ' and θ' represent the target networks. For policy optimization, DQL employs the most basic form of Offline RL optimization, which involves training the policy to maximize the Q-value while imitating behavior policies, using a weighting factor α to balance the influence of both aspects:

$$\mathcal{L}(\theta) = \mathcal{L}_{\text{score}}(\theta) - \alpha \cdot \mathbb{E}_{s \sim \mathcal{D}, a_0 \sim \pi_\theta} [Q_\phi(s, a_0)], \quad (33)$$

where $\mathcal{L}_{\text{score}}(\theta)$ is the score matching loss used for diffusion model training. As the scale of the Q-value function varies in different offline datasets, to normalize it, DQL sets $\alpha = \frac{\eta}{\mathbb{E}_{(s, a) \sim \mathcal{D}} [Q_\phi(s, a)]}$ and tunes η for loss term balance. The Q_ϕ in the denominator is only for normalization and not differentiated over.

Efficient Diffusion Policy. [22] Efficient Diffusion Policy (EDP) aims to address the significant computational overhead caused by iterative sampling and gradient computation during the training of the DQL. Compared to DQL, EDP proposes using DPM-Solver instead of DDPM to reduce the number of sampling steps. Then, EDP introduces an *action approximation* technique, where during policy optimization, one-step denoising is performed on the perturbed action a_t to approximate a_0 . For the process using a data prediction model x_θ and a noise prediction model ϵ_θ separately, the following two equations can express the technique:

$$a_0 \approx x_\theta(a_t, t) \quad (34)$$

$$a_0 \approx \frac{a_t - \sigma_t \epsilon_\theta(a_t, t)}{\alpha_t}. \quad (35)$$

EDP reduces the sampling steps to 15 (even though DQL has only 5 sampling steps) and performs only one-step denoising during policy optimization, significantly speeding up the model training process and achieving performance close to that of DQL.

Implicit Diffusion Q-Learning. [16] Implicit Diffusion Q-Learning (IDQL) models the policy from the perspective of general *constrained policy search* (CPS), in which the optimal policy is described as a weighted behavior policy:

$$\pi^*(a|s) = \pi_\theta^b(a|s) w(a|s), \text{ s.t. } \int_{\mathcal{A}} w(a|s) da = 1, \forall s, \quad (36)$$

where $\pi_\theta^b(a|s)$ represents the behavior policy learned by the diffusion model from the dataset, and $w(s, a)$ is a weight function. IDQL derives its weight function from the generalized implicit

370 Q-learning:

$$w(\mathbf{a}|\mathbf{s}) = \frac{|f'(Q_\phi(\mathbf{s}, \mathbf{a}) - V^*(\mathbf{s}))|}{|Q_\phi(\mathbf{s}, \mathbf{a}) - V^*(\mathbf{s})|}, \quad (37)$$

371 where f can be any convex function, $f' = \frac{\partial f}{\partial V(\mathbf{s})}$, and

$$V^*(\mathbf{s}) = \arg \min_{V(\mathbf{s})} \mathbb{E}_{\mathbf{a} \sim \pi_\theta^b(\mathbf{a}|\mathbf{s})} [f(Q_\phi(\mathbf{s}, \mathbf{a}) - V(\mathbf{s}))]. \quad (38)$$

372 Therefore, the training of IDQL consists of two independent processes: training the diffusion model
 373 to clone the behavior policy and training the IQL-based weight function $w(\mathbf{a}|\mathbf{s})$. At each inference
 374 step, IDQL samples a set of candidate actions $\{\mathbf{a}_0\}$, computes the weights $\{w(\mathbf{s}, \mathbf{a}_0)\}$, and then
 375 selects the action to be executed as a categorical from $\{w(\mathbf{s}, \mathbf{a}_0)\}$.

376 **DiffusionBC.** [40] DiffusionBC constructs an observation-to-action diffusion model for imitating
 377 stochastic and multimodal human demonstrations. The basic version of DiffusionBC applies diffusion
 378 generation directly as a diffusion policy $\pi(\mathbf{a}_0|\mathbf{s}, \mathbf{a}_t, t)$ with noisy action $\mathbf{a}_t \in \mathbb{R}^{|\mathbf{a}|}$, denoising timestep
 379 t and observation \mathbf{s} (possibly with a history) input. To better select intra-distributional actions to
 380 mimic human behavior, DiffusionBC proposed the **Diffusion-X Sampling** trick, which encourages
 381 higher likelihood actions during sampling. For diffusion-X sampling, the sampling process first
 382 runs normal T denoising timesteps, and timesteps is fixed to $t = 1$, then extra denoising iterations
 383 continue to run for M timesteps toward higher-likelihood regions.

384 **DiffusionPolicy.** [4] Similar to DiffusionBC, Diffusion Policy also uses a diffusion model to directly
 385 approximate the conditional distribution $p(\mathbf{a}|\mathbf{s})$, but uses two key design choices: (1) **Closed-loop**
 386 **Action-chunking Prediction:** Diffusion Policy generates sequences of actions per prediction rather
 387 than single action to encourage temporal consistency and smoothness in long-term planning to better
 388 fit multimodal distributions. At time step t , the policy takes the latest T_s (the observation horizon)
 389 steps of observation data \mathbf{s}_t as input and predicts H steps of actions, of which T_a (the action prediction
 390 horizon) steps of actions are executed on the robot without re-planning. (2) **Network Architecture**
 391 **Options:** Diffusion Policy adopts the traditional 1D-Unet [21] and DiT [41] to new CNN-based Unet
 392 and time-series diffusion transformer network architectures. CNN-based Diffusion Policy conditions
 393 the action generation process on observation \mathbf{s} with Feature-wise Linear Modulation (FiLM) [42] and
 394 Transformer-based Diffusion Policy fuses state \mathbf{s} and action \mathbf{a} features via cross attention to jointly
 395 predict $\epsilon_\theta(o, a_k, k)$, where k is sinusoidal embedding for diffusion iteration. The Diffusion Policy
 396 has demonstrated excellent performance and high stability in multiple simulation environments and
 397 real-world tasks for imitation learning and is a widely used baseline for embodied AI.

398 G.3 Diffusion Data Synthesizers.

399 **SynthER.** [36] SynthER uses the diffusion model to generate one-step transitions $(\mathbf{s}, \mathbf{a}, r, d, \mathbf{s}')$.
 400 Trained on an offline dataset, SynthER then upsamples it to a larger dataset (in D4RL, SynthER
 401 upsamples each dataset to 5M transitions), which helps other offline RL algorithms to optimize the
 402 agent policy.

403 H Limitations, Challenges, and Future Directions

404 **Limitations.** Although the modular structure and pipeline design of CleanDiffuser greatly simplify
 405 the implementation difficulty for researchers deploying DMs, the inherent complexity of the principles
 406 and improvements of DMs still requires a considerable amount of time to deeply understand each type
 407 of module. We hope to alleviate this issue and better facilitate collaboration through comprehensive
 408 configuration files and documentation, as well as active maintenance and updates. Additionally,
 409 When dealing with certain specific issues, CleanDiffuser may require tailored adjustments and
 410 optimizations. For instance, the current version of CleanDiffuser does not directly support
 411 discrete or hybrid action space tasks, which may be mitigated through techniques such as action
 412 representation [28] or using categorical diffusion models [6].

413 Based on experimental analyses of CleanDiffuser, we have identified several promising areas for
414 further research as follows:

415 **Unleashing the potential of diffusion planners.** Analogous to the classification of RL algorithms,
416 as diffusion planners can imaginatively generate interactive trajectories, they should be categorized
417 under model-based RL (MBRL). In MBRL, there are various ways to utilize learned dynamic
418 models, including planning to search for the optimal action [13, 17], optimizing policies using
419 rollout trajectories [12], and even combining these two approaches [15, 14]. Currently, diffusion
420 planners are limited to the first paradigm, and due to their sensitivity to guidance and lack of safety
421 constraints, they are prone to OOD plans [7], falling short in performance compared to other offline
422 MBRL algorithms. Future research can explore new paradigms for diffusion planners, attempting
423 diverse ways to utilize generated trajectories or integrating safety constraints to enhance the fidelity
424 of generated trajectories, thereby unleashing the full potential of diffusion planners.

425 **Exploring the reasons behind sampling degradation.** In ??, we discuss an anomaly known as
426 *sampling degradation*, where the algorithm’s performance decreases as the number of sampling steps
427 increases. This anomaly has been identified in previous works [22, 3] and remains an open question.
428 Theoretically, more sampling steps should result in a more accurate SDE/ODE solution, ultimately
429 producing higher-fidelity samples. This naturally prompts a trade-off exploration between sampling
430 steps and performance during implementation. However, in experiments, increasing sampling steps
431 in certain tasks does not improve performance and can even lead to a decrease. Future research can
432 systematically investigate this anomaly to provide optimal recommendations for selecting sampling
433 steps.

434 **Understanding the impact of SDE and ODE.** In our experiments, we observe consistent differences
435 in SDE solvers and ODE solvers on algorithm performance, tendency to sampling degradation, and
436 sensitivity to guidance. While there is existing research on the impact of SDE and ODE in computer
437 vision [39, 35], there is still a gap in research within the decision-making domain. Future research
438 can fill this gap and explore the implications of SDE and ODE solvers in decision-making tasks.

439 **Accelerating Diffusion Model Sampling.** Due to the denoising process involved in iterative sampling,
440 DMs face the issue of slow sampling speeds when used for decision-making. This poses significant
441 challenges in scenarios such as real-time robot control or game AI. DiffuserLite [7] is a diffusion
442 planner method that addresses this issue by modeling the diffusion process through a plan refinement
443 process for coarse-to-fine-grained trajectory generation and further accelerates the sampling speed
444 using rectified flow. Further speeding up the sampling speed of various roles of DMs remains a
445 promising research direction.

446 I Potential Social Impact

447 CleanDiffuser fills a critical gap in the current landscape by providing a unified and modularized
448 framework that empowers researchers and practitioners to explore new frontiers. This will accelerate
449 the development and deployment of diffusion-based decision-making applications, such as various
450 robotics research and products. However, CleanDiffuser may also be used in military weapon
451 development.

452 J License

453 Our codebase is released under Apache License 2.0.

References

- [1] Anurag Ajay, Yilun Du, Abhi Gupta, Joshua B. Tenenbaum, Tommi S. Jaakkola, and Pulkit Agrawal. Is conditional generative modeling all you need for decision making? In The Eleventh International Conference on Learning Representations, ICLR, 2023.
- [2] Greg Brockman, Vicki Cheung, Ludwig Pettersson, Jonas Schneider, John Schulman, Jie Tang, and Wojciech Zaremba. Openai gym. arXiv preprint arXiv:1606.01540, 2016.
- [3] Huayu Chen, Cheng Lu, Chengyang Ying, Hang Su, and Jun Zhu. Offline reinforcement learning via high-fidelity generative behavior modeling. In The Eleventh International Conference on Learning Representations, ICLR, 2023.
- [4] Cheng Chi, Siyuan Feng, Yilun Du, Zhenjia Xu, Eric Cousineau, Benjamin Burchfiel, and Shuran Song. Diffusion policy: Visuomotor policy learning via action diffusion. In Proceedings of Robotics: Science and Systems, RSS, 2023.
- [5] Prafulla Dhariwal and Alexander Quinn Nichol. Diffusion models beat GANs on image synthesis. In Advances in Neural Information Processing Systems, NIPS, 2021.
- [6] Sander Dieleman, Laurent Sartran, Arman Roshannai, Nikolay Savinov, Yaroslav Ganin, Pierre H Richemond, Arnaud Doucet, Robin Strudel, Chris Dyer, Conor Durkan, et al. Continuous diffusion for categorical data. arXiv preprint arXiv:2211.15089, 2022.
- [7] Zibin Dong, Jianye Hao, Yifu Yuan, Fei Ni, Yitian Wang, Pengyi Li, and Yan Zheng. Diffuserlite: Towards real-time diffusion planning. arXiv preprint arXiv:2401.15443, 2024.
- [8] Pete Florence, Corey Lynch, Andy Zeng, Oscar A Ramirez, Ayzaan Wahid, Laura Downs, Adrian Wong, Johnny Lee, Igor Mordatch, and Jonathan Tompson. Implicit behavioral cloning. In Conference on Robot Learning, CoRL, 2022.
- [9] Justin Fu, Aviral Kumar, Ofir Nachum, George Tucker, and Sergey Levine. D4rl: Datasets for deep data-driven reinforcement learning. arXiv preprint arXiv:2004.07219, 2020.
- [10] Scott Fujimoto and Shixiang Gu. A minimalist approach to offline reinforcement learning. In Advances in Neural Information Processing Systems, NIPS, 2021.
- [11] Abhishek Gupta, Vikash Kumar, Corey Lynch, Sergey Levine, and Karol Hausman. Relay policy learning: Solving long-horizon tasks via imitation and reinforcement learning. In Proceedings of the Conference on Robot Learning, CoRL, 2020.
- [12] Danijar Hafner, Timothy Lillicrap, Jimmy Ba, and Mohammad Norouzi. Dream to control: Learning behaviors by latent imagination. In International Conference on Learning Representations, ICLR, 2020.
- [13] Danijar Hafner, Timothy Lillicrap, Ian Fischer, Ruben Villegas, David Ha, Honglak Lee, and James Davidson. Learning latent dynamics for planning from pixels. In Proceedings of the 36th International Conference on Machine Learning, ICML, 2019.
- [14] Nicklas Hansen, Hao Su, and Xiaolong Wang. TD-MPC2: Scalable, robust world models for continuous control. In The Twelfth International Conference on Learning Representations, ICLR, 2024.
- [15] Nicklas A Hansen, Hao Su, and Xiaolong Wang. Temporal difference learning for model predictive control. In Proceedings of the 39th International Conference on Machine Learning, ICML, 2022.
- [16] Philippe Hansen-Estruch, Ilya Kostrikov, Michael Janner, Jakub Grudzien Kuba, and Sergey Levine. Idql: Implicit q-learning as an actor-critic method with diffusion policies. arXiv preprint arXiv:2304.10573, 2023.

- [17] Xiaotian Hao, Jianye Hao, Chenjun Xiao, Kai Li, Dong Li, and Yan Zheng. Multiagent gumbel muzero: Efficient planning in combinatorial action spaces. Proceedings of the AAAI Conference on Artificial Intelligence, AAAI, 2024.
- [18] Longxiang He, Li Shen, Linrui Zhang, Junbo Tan, and Xueqian Wang. Diffcps: Diffusion model based constrained policy search for offline reinforcement learning. arXiv preprint arXiv:2310.05333, 2024.
- [19] Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. In Advances in Neural Information Processing Systems, NIPS, 2020.
- [20] Jonathan Ho and Tim Salimans. Classifier-free diffusion guidance. In NeurIPS 2021 Workshop on Deep Generative Models and Downstream Applications, 2021.
- [21] Michael Janner, Yilun Du, Joshua Tenenbaum, and Sergey Levine. Planning with diffusion for flexible behavior synthesis. In Proceedings of the 39th International Conference on Machine Learning, ICML, 2022.
- [22] Bingyi Kang, Xiao Ma, Chao Du, Tianyu Pang, and Shuicheng Yan. Efficient diffusion policies for offline reinforcement learning. Advances in Neural Information Processing Systems, NIPS, 36, 2024.
- [23] Tero Karras, Miika Aittala, Timo Aila, and Samuli Laine. Elucidating the design space of diffusion-based generative models. In Alice H. Oh, Alekh Agarwal, Danielle Belgrave, and Kyunghyun Cho, editors, Advances in Neural Information Processing Systems, NIPS, 2022.
- [24] Diederik P Kingma, Tim Salimans, Ben Poole, and Jonathan Ho. Variational diffusion models. In Advances in Neural Information Processing Systems, NIPS, 2021.
- [25] P.E. Kloeden and E. Platen. Numerical Solution of Stochastic Differential Equations. Stochastic Modelling and Applied Probability. Springer Berlin Heidelberg, 2011.
- [26] Ilya Kostrikov, Ashvin Nair, and Sergey Levine. Offline reinforcement learning with implicit q-learning. In International Conference on Learning Representations, ICLR, 2022.
- [27] Aviral Kumar, Aurick Zhou, George Tucker, and Sergey Levine. Conservative q-learning for offline reinforcement learning. In Advances in Neural Information Processing Systems, NIPS, 2020.
- [28] Boyan Li, Hongyao Tang, Yan Zheng, Jianye Hao, Pengyi Li, Zhen Wang, Zhaopeng Meng, and Li Wang. Hyar: Addressing discrete-continuous action reinforcement learning via hybrid action representation. arXiv preprint arXiv:2109.05490, 2021.
- [29] Wenhao Li, Xiangfeng Wang, Bo Jin, and Hongyuan Zha. Hierarchical diffusion for offline decision making. In Proceedings of the 40th International Conference on Machine Learning, ICML, 2023.
- [30] Eric Liang, Richard Liaw, Robert Nishihara, Philipp Moritz, Roy Fox, Ken Goldberg, Joseph Gonzalez, Michael Jordan, and Ion Stoica. RLlib: Abstractions for distributed reinforcement learning. In Proceedings of the 35th International Conference on Machine Learning, ICML, 2018.
- [31] Zhixuan Liang, Yao Mu, Mingyu Ding, Fei Ni, Masayoshi Tomizuka, and Ping Luo. Adapt-diffuser: Diffusion models as adaptive self-evolving planners. In International Conference on Machine Learning, ICML, 2023.
- [32] Haohe Liu, Zehua Chen, Yi Yuan, Xinhao Mei, Xubo Liu, Danilo Mandic, Wenwu Wang, and Mark D Plumbley. AudioLDM: Text-to-audio generation with latent diffusion models. In Proceedings of the 40th International Conference on Machine Learning, ICML, 2023.

- [33] Xingchao Liu, Chengyue Gong, and qiang liu. Flow straight and fast: Learning to generate and transfer data with rectified flow. In The Eleventh International Conference on Learning Representations, ICLR, 2023.
- [34] Cheng Lu, Yuhao Zhou, Fan Bao, Jianfei Chen, Chongxuan Li, and Jun Zhu. DPM-solver: A fast ODE solver for diffusion probabilistic model sampling in around 10 steps. In Advances in Neural Information Processing Systems, NIPS, 2022.
- [35] Cheng Lu, Yuhao Zhou, Fan Bao, Jianfei Chen, Chongxuan Li, and Jun Zhu. Dpm-solver++: Fast solver for guided sampling of diffusion probabilistic models. arXiv preprint arXiv:2211.01095, 2023.
- [36] Cong Lu, Philip Ball, Yee Whye Teh, and Jack Parker-Holder. Synthetic experience replay. Advances in Neural Information Processing Systems, NIPS, 36, 2024.
- [37] Ajay Mandlekar, Danfei Xu, Josiah Wong, Soroush Nasiriany, Chen Wang, Rohun Kulkarni, Li Fei-Fei, Silvio Savarese, Yuke Zhu, and Roberto Martín-Martín. What matters in learning from offline human demonstrations for robot manipulation. In 5th Annual Conference on Robot Learning, CoRL, 2021.
- [38] Alexander Quinn Nichol and Prafulla Dhariwal. Improved denoising diffusion probabilistic models. In Proceedings of the 38th International Conference on Machine Learning, ICML, 2021.
- [39] Shen Nie, Hanzhong Allan Guo, Cheng Lu, Yuhao Zhou, Chenyu Zheng, and Chongxuan Li. The blessing of randomness: SDE beats ODE in general diffusion-based image editing. In The Twelfth International Conference on Learning Representations, ICLR, 2024.
- [40] Tim Pearce, Tabish Rashid, Anssi Kanervisto, Dave Bignell, Mingfei Sun, Raluca Georgescu, Sergio Valcarcel Macua, Shan Zheng Tan, Ida Momennejad, Katja Hofmann, and Sam Devlin. Imitating human behaviour with diffusion models. In The Eleventh International Conference on Learning Representations, ICLR, 2023.
- [41] William Peebles and Saining Xie. Scalable diffusion models with transformers. In Proceedings of the IEEE/CVF International Conference on Computer Vision, ICCV, 2023.
- [42] Ethan Perez, Florian Strub, Harm De Vries, Vincent Dumoulin, and Aaron Courville. Film: Visual reasoning with a general conditioning layer. In Proceedings of the AAAI conference on artificial intelligence, AAAI, 2018.
- [43] Antonin Raffin, Ashley Hill, Adam Gleave, Anssi Kanervisto, Maximilian Ernestus, and Noah Dormann. Stable-baselines3: Reliable reinforcement learning implementations. Journal of Machine Learning Research, 2021.
- [44] Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-resolution image synthesis with latent diffusion models. In Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, CVPR, 2022.
- [45] Nataniel Ruiz, Yuanzhen Li, Varun Jampani, Yael Pritch, Michael Rubinstein, and Kfir Aberman. Dreambooth: Fine tuning text-to-image diffusion models for subject-driven generation. In Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, CVPR, 2023.
- [46] Jiaming Song, Chenlin Meng, and Stefano Ermon. Denoising diffusion implicit models. In International Conference on Learning Representations, ICLR, 2021.
- [47] Yang Song, Jascha Sohl-Dickstein, Diederik P Kingma, Abhishek Kumar, Stefano Ermon, and Ben Poole. Score-based generative modeling through stochastic differential equations. In International Conference on Learning Representations, ICLR, 2021.

- 587 [48] Denis Tarasov, Alexander Nikulin, Dmitry Akimov, Vladislav Kurenkov, and Sergey Kolesnikov.
588 CORL: Research-oriented deep offline reinforcement learning library. In 3rd Offline RL
589 Workshop: Offline RL as a "Launchpad", 2022.
- 590 [49] Patrick von Platen, Suraj Patil, Anton Lozhkov, Pedro Cuenca, Nathan Lambert, Kashif Rasul,
591 Mishig Davaadorj, Dhruv Nair, Sayak Paul, William Berman, Yiyi Xu, Steven Liu, and Thomas
592 Wolf. Diffusers: State-of-the-art diffusion models. [https://github.com/huggingface/](https://github.com/huggingface/diffusers)
593 [diffusers](https://github.com/huggingface/diffusers), 2022.
- 594 [50] Zhendong Wang, Jonathan J Hunt, and Mingyuan Zhou. Diffusion policies as an expressive
595 policy class for offline reinforcement learning. In The Eleventh International Conference on
596 Learning Representations, ICLR, 2023.
- 597 [51] Lvmin Zhang, Anyi Rao, and Maneesh Agrawala. Adding conditional control to text-to-image
598 diffusion models. In Proceedings of the IEEE/CVF International Conference on Computer
599 Vision, ICCV, 2023.
- 600 [52] Tony Z. Zhao, Vikash Kumar, Sergey Levine, and Chelsea Finn. Learning Fine-Grained
601 Bimanual Manipulation with Low-Cost Hardware. In Proceedings of Robotics: Science and
602 Systems, RSS, 2023.