

MODULAR DIFFUSION POLICY TRAINING: DECOUPLING AND RECOMBINING GUIDANCE AND DIFFUSION FOR OFFLINE RL

Anonymous authors

Paper under double-blind review

ABSTRACT

In classifier-free diffusion(CFD), the diffusion model and its guidance are typically learned jointly and applied jointly in the inference stage. Before the guidance has converged, it provides unstable or even misleading gradients, which leads to inefficiency and instability during the early stage of training. Such strict coupling not only leads to self-enforcing variance and biased errors but also prevents the guidance module from being reused across different diffusion models. We propose Guidance-First Diffusion Training (GFDT), which pretrains and freezes the guidance model before diffusion policy learning. GFDT reduces peak memory and computation by 38.1%, decreases diffusion training by 65.6% and 27.66%, and achieves up to 12% and 54% performance improvements on offline RL benchmarks. Beyond efficiency, we uncover a strong plug-and-play property: replacing the guidance module only at inference time can substantially improve stability. For example, swapping seeds reduces Median variance by 14% and 21%, and cross-algorithm swaps (e.g., Implicit Q-Learning (IDQL) guidance for Diffusion Q-Learning (DQL) policies) perform comparably to the stronger of the two, despite never being co-trained. Our theoretical analysis shows that GFDT enables the convergence on an optimal guidance and theoretically proves that it speeds up the training. Also, we proved that plug-and-play remains valid as long as the guidance and the diffusion model are trained with the same data distribution. Limitations arising from dataset mismatch are analyzed in detail, which further underscores the necessity of distributional alignment. This work opens a new line of research by treating diffusion and guidance as modular units that can be recombined, rather than as a monolithic process, suggesting a paradigm that may guide the future development of diffusion-based reinforcement learning.

1 INTRODUCTION

By formulating policy learning as a conditional generative process, **diffusion policies** are capable of modeling complex, multimodal behaviors and have demonstrated strong empirical performance across a wide range of RL and robotic manipulation tasks. A central challenge in diffusion policy learning lies in injecting reward optimization into the denoising process, especially in reinforcement learning, where reward is the main evaluation criterion. One of the most widely used methods is **Classifier-Free Guidance (CFG)**¹, which biases sampling toward high-reward trajectories.

Currently, the guidance module and the diffusion model are trained jointly and tightly coupled in inference. The TWO CHALLENGES that exist are: 1)The role of the guidance module is to steer the diffusion process toward generating actions that lead to higher rewards. However, before the guidance module has converged, it cannot provide effective guidance. Therefore, during the early stages of training, the diffusion model receives little to no meaningful directing from the guidance module, which misguides the diffusion process and degrades its performance Kim et al. (2023). 2) This tightly bound diffusion and guidance model can have common idiosyncratic errors and cause instability(large variance). Finally, the coupled usage prevents reuse across different combinations of network modules.

¹All abbreviations are summarized in the Appendix.L

Existing methods solve CHALLENGE 1 by pretraining the diffusion model and then letting the diffusion model generate samples for the guidance training. In this work, we take the opposite approach: pre-training the guidance module, then freezing it to guide the diffusion model. This method is called Guidance-First Diffusion Training (GFDT), which speeds up early-stage training, reduces memory usage, and outperforms joint training methods. Fundamentally, the guidance module is trained using observations, actions, and rewards from a dataset, data that are independent of the actions generated by the diffusion model. Therefore, the guidance module can be trained separately using supervised learning and then be used to effectively guide the diffusion model. Importantly, this separation does not reduce exploration, as offline reinforcement learning relies entirely on fixed datasets, without any online exploration.

To solve CHALLENGE 2, CFG methods can be modularized: our Double Guidance diffusion architecture addresses this by using a differently seeded version of the guidance network during inference, distinct from the one used in training. Existing research basically assumes identical networks using different seeds are identical, but in practice, the random initialization of networks can result in different complementary inductive biases and idiosyncratic errors. This replacement helps correct inaccurate estimations and significantly reduces variance. Further, we cross-combined components from models trained under different frameworks. Specifically, we took the guidance module from an IDQL model and paired it with a completely separate DQL diffusion model—one that had never been trained in conjunction with the IDQL guidance. The plug-and-play (PAP) module and its reversed model are functional, and the performance is comparable with the stronger of these two, and has a significantly better early training stage, outperforming both the baselines. This result strongly suggests that the relationship between the guidance module and the diffusion model is modular and flexible, and therefore, highlights a promising direction for future research.

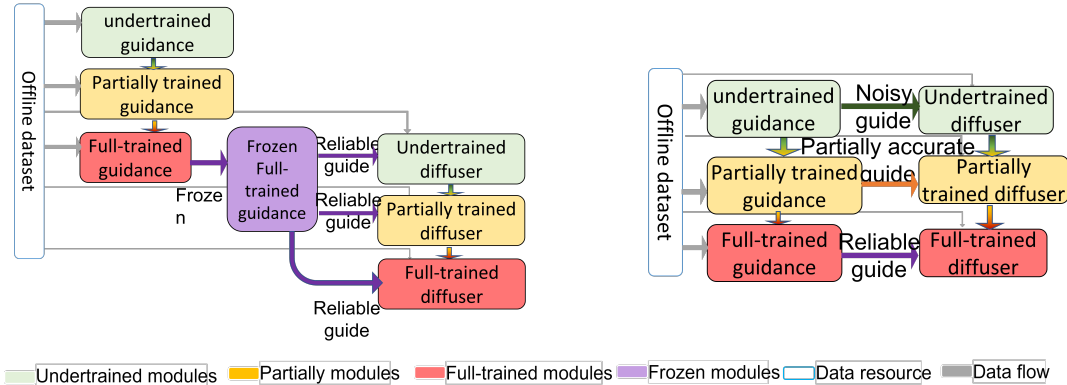


Figure 1: Comparison between the training stage(left) and inference stage(right) of GFDT

The detailed pseudo-codes of the GFDT and the traditional CFG, and a diagram comparison can be found in Appendix C.1.

2 PRELIMINARIES

In offline reinforcement learning (RL), the objective is to learn a policy $\pi_\theta(a|s)$ that maximizes the expected return using only a fixed dataset $D = \{s, a, r, s'\}$, with interaction to the environment only at the inference stage Levine et al. (2020). When diffusion models are used as policies, several offline RL methods—including DQL Janner et al. (2022), Exponential Diffusion Process(EDP) Kang et al. (2023), and IDQL Hansen-Estruch et al. (2023)—adopt a *conditional denoising diffusion process*. The policy is defined implicitly by the reverse denoising procedure: $\hat{a}_0 \sim \pi_\theta(\cdot | s)$. The details of Q-Guided Diffusion are in Appendix B.

Diffusion Q-learning (DQL): The DQL model modified the traditional Q-guided diffusion and applied a variation of the Deep Q network called double Q to improve the accuracy of the Q-value estimation van Hasselt (2010).

$$\mathcal{L}_\pi^{\text{DQL}}(\theta) = -\mathbb{E}_{s \sim \mathcal{D}} [\min (Q_{\phi_1}(s, a^0), Q_{\phi_2}(s, a^0))], \quad \text{where } a^0 \sim \pi_\theta(\cdot | s) \quad (1)$$

Implicit Diffusion Q-learning (IDQL). Unlike DQL, which directly incorporates the reward target into the diffusion training, IDQL adopts a two-stage design. During training, the diffusion model is updated purely via behavior cloning from the dataset, without any explicit Q-guidance. In parallel, a separate Q-value network is learned to estimate $Q(s, a)$. At inference time, the diffusion model offers candidate actions, and the guidance module evaluates them with the Q-network through a one-hot encoding scheme, and selects the action with the highest estimated Q-value. Compared to traditional Q-guided diffusion, IDQL thus applies the guidance entirely to the inference stage rather than the training stage. Although IDQL does not leverage the Q-value estimator to directly guide diffusion during training, the implementation still places the Q-network update within the same training loop as the diffusion model. This coupling is not theoretically necessary, and the consequence is a substantially higher computational and memory usage. There is a detailed analysis on the different between IDQL and GFDT in Appendix C.2.

Efficient Diffusion Policy (EDP) introduces a key modification on traditional Q-guided diffusion: *one-step denoising*. Instead of running the full reverse chain to generate actions, EDP corrupts a dataset action a^0 to a^k in one step, and then one-step-reconstructs an approximate clean action \hat{a}^0 using the denoise backward pass as Equation 6 in Appendix B. This *action approximation* replaces expensive sampling with a lightweight inference step, making EDP orders of magnitude faster to train. We only provide a brief summary here, and defer further details to Appendix D on the details of the one-step guidance application.

These three methods are the baselines of our methods. Currently, our modifications are typically only applied to TD-based methods and cannot be applied to Trajectory-Based methods(e.g., Diffuser Janner et al. (2022); Ajay et al. (2023)) as they use return annotations for full sequences $\tau = (s_0, a_0, \dots, s_T)$. However, experiments showed low convergence rates and suboptimal precisions in their Q-value predictions, and therefore, Guidance-First Diffusion Training showed little improvement. The reason for this observation is supposed to be: when the trajectory is long, the return may be noisy due to stochastic behavior policies, so reward attribution over entire sequences can be ambiguous. (It is hard to infer which action caused the reward change.) In contrast, TD-based methods that work on (s, a) pairs avoid this issue and support more granular, local learning signals. See Appendix E for more explanations.

3 METHODOLOGY

Having introduced the GFDT method and modular module composition in the introduction, we now extend the discussion with a mathematical model to justify the approach. First, we explain why a pretrained, converged model can provide more accurate guidance to the diffusion model, better than an unconverged model that is trained jointly with it, thereby accelerating the training process. Second, we prove that as long as the guidance model can provide a direction that leads to reward improvement and the step size is small enough, it is sufficient for guiding the diffusion model, even if the guidance model was not co-trained with the diffusion model.

Theoretical Motivation. We build upon the theoretical framework established by Theorem Fujimoto et al. (2019a), which introduces Batch-Constrained Q-Learning (BCQL) as a value-based method with provable convergence guarantees under offline settings. *Given a deterministic MDP and coherent batch \mathcal{B} , along with standard Robbins-Monro convergence conditions on the learning rate, BCQL converges to $Q^{\pi^{\mathcal{B}}}(s, a)$, where $\pi^*(s) = \arg \max_{a \text{ s.t. } (s,a) \in \mathcal{B}} Q^{\pi^{\mathcal{B}}}(s, a)$. This policy is guaranteed to match or outperform any behavioral policy contained in the dataset.*

Batch-Constrained Guidance via Diffusion. Inspired by this result, we adopt a *batch-constrained guidance* framework for training diffusion policies. Our key design choice is to pretrain a guidance policy (e.g., a value function $Q(s, a)$ or behavioral prior) purely on the offline dataset, then use it to guide the sampling of a diffusion policy. Theorem 1 guarantees, the training on the offline dataset converges to an optimal reward estimation policy $Q(s,a)$.

This optimized $Q(s,a)$ is then added to the gradient update of the training process of the diffusion:

$$\nabla_{\theta} \mathcal{L}_{\text{total}} = \nabla_{\theta} \mathcal{L}_{\text{BC}} + \eta \nabla_{\theta} Q(s, a), \quad (2)$$

where \mathcal{L}_{BC} is a behavior cloning loss that regularizes the learned policy to stay close to the empirical distribution of \mathcal{B} , and η controls the strength of the guidance. Both of these two parts are learned

based on the same dataset \mathcal{B} , so term 1 offers regularization to avoid leaving the support of the offline dataset and term 2 encourages reward optimization.

To further ensure batch-constrained optimization behavior during action generation(inference step), we modify the reverse diffusion process by injecting Qvalue gradients of the pretrained module:

$$a_{t-1} \leftarrow f_{\text{diffusion}}(a_t, \hat{a}_0, t) + \lambda \frac{\nabla_a Q(s, \hat{a}_t)}{\|\nabla_a Q(s, \hat{a}_t)\| + \varepsilon} + \sqrt{2\tau_t} \xi_t,$$

where λ is a guidance coefficient. This operation encourages the final action a_0 to lie in high-reward regions. Importantly, because the value function $Q(s, a)$ is trained on the offline dataset \mathcal{B} , and the diffusion model itself models the same data distribution. Intuitively, in every inference step, the diffusion model generates an action, and the Q module shifts the gradient in a magnitude that is smaller than or equal to λ . So, the generated action must be inside the region of the dataset \mathcal{B} . Adding the small gradient shift, the entire action is still inside or close to the dataset \mathcal{B} 's coverage. Thus, the value-aware perturbation can be interpreted as a form of approximate, soft batch-constrained policy improvement. For a more detailed explanation, please check Appendix G

On the contrary, before sufficient training, the inaccurate gradient will misshape the diffusion distribution and prolong the training process by adding incorrect information to the optimization, which may or may not be corrected in the later training.

Admittedly, our approach does not strictly satisfy all assumptions of BCQL: the network may converge not to the true Bellman optimum but to a local minimum, an inherent limitation of non-convex optimization that cannot be fully avoided. However, we inherit its core principle: The guidance can be trained to an optimal value estimator, and any value-based guidance must be constrained to the data distribution to guarantee that the gradient guidance is reliable. This prove is also consistent with the experiment results (Section 4): in distribution guidance improves the performance of the CFG. Validated by inverse reasoning, the application scope of the GFDT method in Section 5 showed that out-of-distribution guidance cannot guarantee efficient guidance, and in practice, often jeopardizes the performance.

Theorem 2: Joint training of the guidance module and the diffusion policy is not required. A pretrained Q-network Q_ϕ can be modularized and directly applied to guide the diffusion model, as long as Q_ϕ accurately estimates state-action values within the dataset support \mathcal{B} .

We consider the perturbed diffusion sampling process after the action generation:

$$a_t = a_t + \lambda \cdot \frac{\nabla_a Q_\phi(a_t)}{\|\nabla_a Q_\phi(a_t)\|} + \sqrt{2\tau_t} \xi_t, \quad (3)$$

where λ is a small step size and τ_t controls the annealed noise. The gradient $\nabla_a Q_\phi$ provides a reward-sensitive vector field over the action space.

Directional Alignment and Convergence. As long as the value gradient is directionally positive correlation the diffusion model's score function(i.e., $\cos \theta = \angle(\nabla Q_\phi, \nabla \log p(a)) > 0$), then the expected change in Q_ϕ over one step is

$$\mathbb{E}[Q_\phi(a_{t+1}) - Q_\phi(a_t)] \approx \lambda \cdot \cos \theta_t \cdot \|\nabla Q_\phi(a_t)\| + \mathcal{O}(\sqrt{\tau_t}). \quad (4)$$

Moreover, since the noise is Gaussian, we have: $\lim_{T \rightarrow \infty} \frac{1}{T} \sum_{t=1}^T \xi_t = 0$ (in expectation), meaning the cumulative effect of noise diminishes over time, and the value gradient dominates in expectation. The action will be guided towards a rewarding-increasing direction gradually.

A pretrained Q-function can reliably steer the diffusion process toward high-reward regions without destabilizing sampling, as long as the step size is small and the guidance remains positively beneficial and smooth. This justifies the plug-and-play design of our framework and aligns with the batch-constrained principle in offline RL.

Why Pretraining Helps: Reducing Optimization Burden For a guidance model Q_ϕ trained to ϵ -accuracy on dataset \mathcal{B} , the gradient estimation error satisfies:

$$\mathbb{E}\|\nabla_a Q_\phi - \nabla_a Q^*\| \leq \frac{L}{\sqrt{N}} + \epsilon \quad (5)$$

$\nabla_a Q_\phi$	Pretrained Guidance Gradient	ϵ	Final training error of Q_ϕ
L	Lipschitz constant of the Q-function class	N	Number of samples in dataset \mathcal{B}
$\nabla_a Q^*$	True gradient of the optimal Q-function		

Heuristically, although in reinforcement learning practice the Lipschitz constant L is rarely specified, in related areas (e.g., WGAN with spectral normalization Miyato et al. (2018)) one often enforces $L = 1$. Under this convention, Eq.5 implies that achieving a guidance accuracy of $\mathbb{E}\|\nabla_a Q_\phi - \nabla_a Q^*\| = 0.01$ would require about $N = 1/\mathbb{E}\|\nabla_a Q_\phi - \nabla_a Q^*\|^2 = 10^4$ effective samples (equivalently, $\sim 10^4$ gradient updates of the Q network). In more parameter-sensitive scenarios, pushing the accuracy further to $\mathbb{E}\|\nabla_a Q_\phi - \nabla_a Q^*\| = 0.001$ would require as many as $N = 10^6$ effective samples (i.e., 10^6 gradient updates). This illustrates the significant optimization burden, and shows how our method helps reduce otherwise wasted early training of the diffusion model. Pretraining is therefore essential to ensure that guidance has semantic meaning before being applied in generation.

4 EXPERIMENTS

We evaluate our method on standard tasks from the PyBullet D4RL benchmark Fu et al. (2021); Foundation (2024). To ensure fair comparison, we re-train three seeds of three representative diffusion-based offline RL algorithms—EDP, DQL, and IDQL(module interchange)—based on publicly available code from Wang et al. (2024). These methods serve as illustrative examples of our theory, which can be easily extended to more diffusion methods. For more details about the environment setting please check Appendix F. Very importantly, the detail comparison between our model and the state-of-the-art models is shown in Table 12 in Appendix H.

4.1 PERFORMANCE OF GFDT

All experiments were conducted using PyTorch on a high-performance computing (HPC) cluster. We adhered to any assigned framework for our work because the focus of this paper is not restricted to any specific computational framework. The comparison between our model and other RL models is in Table 12 of Appendix H

Table 1: Mujoco and Antmaze Results DQL

Env	Baseline	GFDT	Double_GFDT	D_baseline	GAI	Behavior Clone	Unfreeze
HCEX	88.21±0.28	90.43±0.26	90.51±0.29	88.52±0.30	68.49±1.55	85.82±2.24	89.50±0.23
HCME	88.28±0.53	90.18±0.38	89.77±0.95	84.76±1.57	89.73±0.43	85.53±1.34	89.76±0.61
HCMR	67.38±0.34	67.93±0.33	68.23±0.39	68.31±0.27	67.76±0.52	55.87±5.23	67.43±0.37
HCMV	59.88±5.40	66.99±0.60	68.95±0.46	53.49±3.28	53.36±1.43	54.57±4.27	55.04±2.27
HOEX	166.76±0.58	172.90±0.56	172.71±0.60	163.81±1.49	162.68±1.21	165.11±0.67	155.15±44.68
HOME	168.00±1.18	172.31±1.09	167.01±0.86	160.93±1.56	166.28±1.03	166.88±1.33	159.10±31.75
HOMR	151.59±0.57	119.91±11.62	152.98±0.75	156.26±0.93	121.97±49.30	113.30±35.20	150.50±3.44
HOMV	143.92±1.05	147.06±1.09	145.36±1.38	132.60±5.97	71.09±35.69	118.08±15.41	144.59±1.09
WAEX	117.25±0.46	120.25±0.60	119.69±0.55	115.47±0.58	119.91±0.52	108.05±28.27	117.92±0.47
WAME	117.73±0.35	117.63±0.34	112.87±22.44	113.61±2.29	115.08±1.32	117.45±0.68	118.50±0.73
WAMR	92.96±0.60	95.62±4.78	93.28±0.55	43.54±22.52	80.32±19.52	87.14±1.87	95.62±4.78
WAMV	87.65±0.40	87.89±0.27	87.94±0.33	76.72±1.34	77.33±13.13	82.28±0.70	87.89±0.43
Env	Base	GFDT	Double_GFDT	D_baseline	GAI		
large-diverse	63.33±6.94	90.67±5.39	84.67±6.00	84.67±6.00	84.67±6.00	86.67±5.83	
large-play	90.00±5.48	88.00±5.70	89.33±5.56	88.00±5.70	88.00±5.70	88.67±5.63	
medium-diverse	93.33±4.99	97.33±4.01	89.33±5.56	87.33±5.77	87.33±5.77	94.00±4.87	
medium-play	67.33±6.85	91.33±5.30	92.00±5.21	76.00±6.54	88.67±5.63		

From the Table 2 and Table 1, in all environments and algorithms, GFDT outperformed baseline models except 1 has a lower performance. The early performance of the GFDT is significantly better than the baseline model. The improve of the DQL model is higher than the EDP model. Formally, DQL can be abstracted as $x_t = f(x_{t+1}, \text{guidance}_t)$ where each step depends on both the previous state x_t and a step-wise guidance signal, leading to an iterative accumulation of guidance effects. In contrast, EDP follows a single-step mapping $x_0 = f(x_T, \text{guidance})$, where the guidance is injected only once without iterative updates, which naturally limits its overall influence.

Table 2: Mujoco Results and Antmaze Results EDP

Env	Baseline	GFDT	Double_GFDT	D_baseline	GAI	Behavior Clone	Unfreeze
HCEX	88.21±0.28	90.43±0.26	90.51±0.29	88.52±0.30	68.49±1.55	85.82±2.24	89.50±0.23
HCME	88.28±0.53	90.18±0.38	89.77±0.95	84.76±1.57	89.73±0.43	85.53±1.34	89.76±0.61
HCMR	67.38±0.34	67.93±0.33	68.23±0.39	68.31±0.27	67.76±0.52	55.87±5.23	67.43±0.37
HCMV	59.88±5.40	66.99±0.60	68.95±0.46	53.49±3.28	53.36±1.43	54.57±4.27	55.04±2.27
HOEX	166.76±0.58	172.90±0.56	172.71±0.60	163.81±1.49	162.68±1.21	165.11±0.67	155.15±44.68
HOME	168.00±1.18	172.31±1.09	167.01±0.86	160.93±1.56	166.28±1.03	166.88±1.33	159.10±31.75
HOMR	151.59±0.57	119.91±11.62	152.98±0.75	156.26±0.93	121.97±49.30	113.30±35.20	150.50±3.44
HOMV	143.92±1.05	147.06±1.09	145.36±1.38	132.60±5.97	71.09±35.69	118.08±15.41	144.59±1.09
WAEX	117.25±0.46	120.25±0.60	119.69±0.55	115.47±0.58	119.91±0.52	108.05±28.27	117.92±0.47
WAME	117.73±0.35	117.63±0.34	112.87±22.44	113.61±2.29	115.08±1.32	117.45±0.68	118.50±0.73
WAMR	92.96±0.60	95.62±4.78	93.28±0.55	43.54±22.52	80.32±19.52	87.14±1.87	95.62±4.78
WAMV	87.65±0.40	87.89±0.27	87.94±0.33	76.72±1.34	77.33±13.13	82.28±0.70	87.89±0.43
Env	Base	GFDT	Double_GFDT	D_baseline	GAI		
large-diverse	30.67±6.99	31.33±6.91	34.67±6.90	53.33±7.36	28.67±6.83		
large-play	21.33±6.40	22.67±6.70	22.67±6.59	26.67±6.65	18.67±6.37		
medium-diverse	67.33±9.06	52.00±7.72	42.67±6.86	19.33±6.41	2.00±3.74		
medium-play	73.30±10.34	118.00±10.35	114.67±9.66	110.00±9.96	90.00±9.78		

Table 3: Training gradient steps to get 95% performance(batch size 256)

Env	DQL	GFDT_DQL	EDP	GFDT_EDP	dql_diff-idql_guid	idql_diff-dql
HCEX	7600	2800	36.84%	18000	3600	20.00%
HCME	16400	2800	17.07%	31200	6000	19.23%
HCMR	8000	3200	40.00%	20400	8400	41.18%
HCMV	52000	7600	14.62%	10800	10000	92.59%
HOEX	30800	15200	49.35%	8400	8400	100.00%
HOME	41600	23600	56.73%	24000	16800	70.00%
HOMR	24800	3200	12.90%	40800	52800	129.41%
HOMV	79200	12800	16.16%	50400	76800	152.38%
WAEX	59600	14000	23.49%	10800	7200	66.67%
WAME	96000	72400	75.42%	21600	18000	83.33%
WAMR	26000	9600	36.92%	38400	10800	28.13%
WAMV	10800	3600	33.33%	79200	51600	65.15%
AVG			34.40%			72.34%
			128	256	512	
			Diffusion Params	172,230 (70.0%)	174,534 (38.1%)	172,230 (13.7%)
			Guidance Params	73,986 (30.0%)	283,650 (61.9%)	1,082,370 (86.3%)
			Total Memory	0.94 MB	1.75 MB	4.79 MB

Table 4: Parameter statistics summary for different model widths. (This model has a hidden layer size of 256.) Explanation of this table are in AppendixA

4.2 THE ROLE OF THE GUIDANCE MODULE

To analyze the role of reward guidance and whether the guidance is removable or replaceable, we conducted an ablation study of a series of training sessions. In this experiment, we compare the performance of GFDT, baseline models, an abolition study that removed the reward guidance part

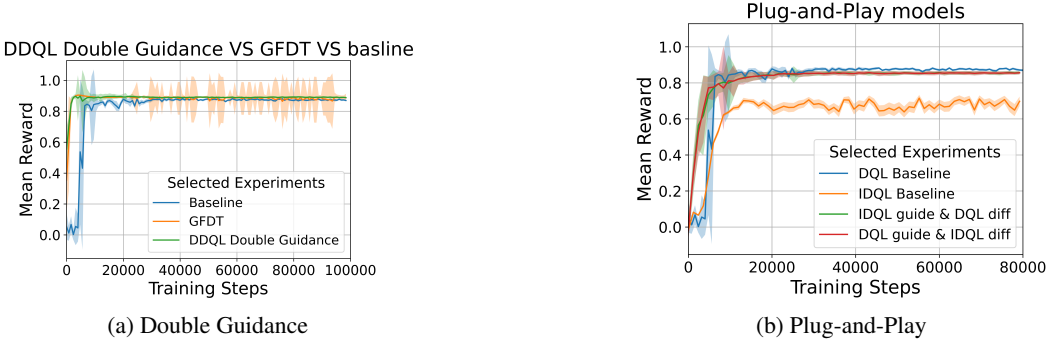


Figure 2: Comparison across Half Cheetah Expert tasks using different algorithms. The legend below applies to all three subplots.

of the baseline training (BC), an abolition study on a behavior clone training with a guidance at the inference stage(GAI), an abolition study that pretrained the guidance module but did not freeze it in diffusion training, two double guidance models as shown in the double guidance section, for the baseline model and the GFDT.

As shown in Figure 2, the pretrained but non-frozen guidance model in diffusion training outperforms the baseline but does not surpass GFDT, confirming that additional training of the guidance model is unnecessary. The slight performance drop is likely caused by overfitting of the guidance model under over-training. Removing the guidance module results (in training) got worse performance and higher variance, underscoring the important role of guidance. Finally, random guidance leads to complete failure of the diffusion model, highlighting that inaccurate guidance signals can be detrimental. This supports our claim that a valid guidance model is critical for successful training.

Admittedly, the increase in training speed is at the cost of foregoing the guidance model first. But the training time of the diffusion model decreased by 65.4% and 27.66% with decreased computational resources. So, the pretraining is very cost-effective. The following sections discuss how to plug and play a pretrained module into the diffusion model, making the modules reusable. The performance of EDP is not as significant because the one-step denoise is more sensitive to conditions.

4.3 DOUBLE GUIDANCE REDUCES VARIANCE

As shown in Figure 2a, the Double Guidance module is defined as replacing the guidance module of a model in the inference stage with a separately trained guidance module that has an identical architecture but differs in its random seed. For example, the guidance module with seed 0 is used to guide the diffusion training (diffusion model still training with random seed 0), but at the inference stage, the guidance module is a different module seeded with 42. The Double Guidance exhibits significantly lower variance, along with all the evaluation results of the checkpoints from training steps 0 to 10000(convergence), than the default models and the baseline model. As shown in the Table.4.3, while the performance of the interchanged is still preserved, the variance has decrease significantly. This is analogous to seeking a second opinion: an independently initialized module reviews the outputs, reducing shared biases and stabilizing the results.

This phenomenon can be explained as the breaking of a self-reinforcing bias loop: when a guidance uses its own (possibly biased) estimates to guide the training of a diffusion model, this behavior may amplify over- or under-estimations. In contrast, using a structurally identical but independently trained guidance model as guidance can mutual correct the bias of each other, similar in spirit to the role of the target network in Double Q-learning. The double guidance model is sufficient for the baseline model; however increases the variance of GFDT because GFDT is train with one guidance all the time and has less experience with a different guidance.

Table 5: Performance Comparison of variance Groups in DQL and EDP

	DQL $\times(10^{-2})$				EDP $\times 10^{-3}$			
	Mean Variance		Median Variance		Mean Variance		Median Variance	
BL	6.49	100.00	4.75	100.00%	4.76	100.00%	1.06	100.00%
D_BL	5.24	80.71	4.11	86.49%	4.64	97.62%	0.85	79.76%
GFDT	6.07	93.46	5.63	118.40%	3.97	83.38%	1.92	180.30%
D_GFDT	8.82	135.93	8.01	168.52%	6.15	129.40%	0.91	85.31%

4.4 PLUG-AND-PLAY MODEL COMPOSITION

Our research reveals that diffusion models exhibit unique plug-and-play compatibility with their guidance modules, following the theoretical proof. We evaluate two hybrid configurations without any additional training: 1) DQL-as-Guidance + IDQL-as-Diffusion, 2) IDQL-as-Guidance + DQL-as-Diffusion. as shown in Figure 3, both combinations: (1) achieve final performance comparable to the DQL baseline, the higher one of these two models, (2) exhibit 34.14% and 40.60% faster initial convergence compared to baselines, and (3) maintain stability despite architectural misplacement.

Using checkpoints (saved every 400 steps), we observe the following properties: First, guidance modules provide effective policy improvement signals regardless of diffusion model architecture. Second, the normalizing step in the training and inference parts is expected to be an important reason

why these two modules are interchangeable, as explained in the Methodology Section Theorem 2. Finally, early-stage advantages suggest shows. When the modules are both under-trained and are not perfect, a differently structured model is not only unharmed, but can cross-correct with the errors of each module. We find that a guidance module—trained independently and even with a completely different architecture—can be directly applied to a diffusion model. This experiment is consistent with our theoretical analysis that, as long as the guidance model provides a reliable estimate of the reward (gradient estimate per se), it can be reused across structurally dissimilar models. Such modularity enables flexible training pipelines and paves the way for reusable, task-agnostic reinforcement learning components. Finally, as we have mentioned in the theoretical part, if the guidance model is not independently trained on the same data distribution, there is no guarantee that the independent module will be beneficial. In the following section, we illustrate that pretraining on out-of-distribution(OOD) data is more likely to be detrimental than beneficial in practice.

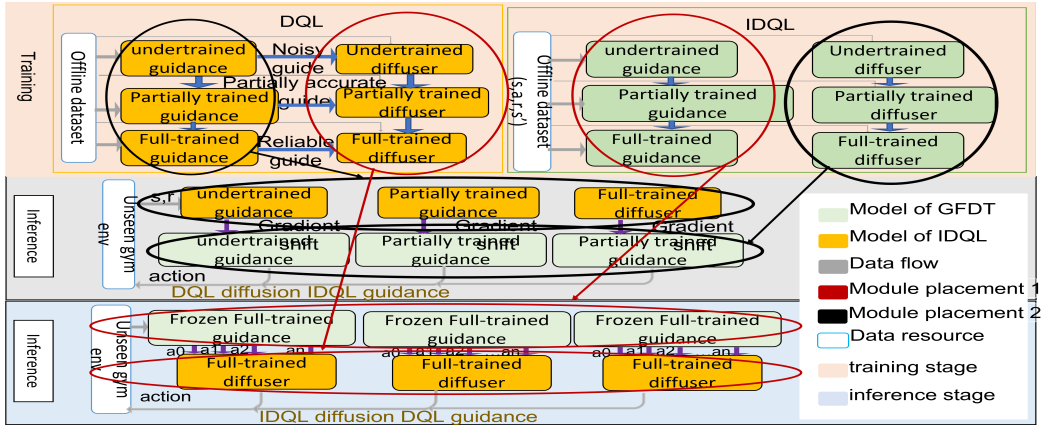


Figure 3: plug-and-play

Table 6: Performance comparison of different algorithms on various datasets

Dataset	DQL	IDQL	DQL-guide-IDQL-guid	IDQL-guide-DQL-guid
halfcheetah-expert	88.21±0.28	71.07±0.02	85.59±6.92	81.64±0.86
halfcheetah-medium-expert	88.28±0.53	72.03±0.01	85.61±0.01	84.77±0.00
halfcheetah-medium-replay	67.38±0.34	55.90±0.00	70.58±0.84	69.79±0.01
halfcheetah-medium	59.88±5.40	52.70±0.00	60.87±0.89	60.87±0.23
hopper-expert	166.76±0.58	160.53±0.00	159.18±4.11	208.67±0.02
hopper-medium-expert	168.00±1.18	128.25±2.82	182.47±0.03	182.68±0.02
hopper-medium-replay	151.59±0.57	55.53±0.84	157.24±0.00	157.24±0.00
hopper-medium	143.92±1.05	65.94±0.26	136.53±0.23	136.53±0.20
walker2d-expert	117.25±0.46	113.78±0.00	116.64±0.14	116.64±0.02
walker2d-medium-expert	117.73±0.35	70.79±0.02	118.11±0.00	118.39±0.00
walker2d-medium-replay	92.96±0.60	71.83±0.05	93.10±0.94	87.35±0.07
walker2d-medium	87.65±0.40	66.78±0.01	87.11±0.07	87.11±0.01
Average	112.47	82.09	112.75	115.97

5 APPLICATION SCOPE AND LIMITATIONS

To illustrate the distribution shift between datasets, We performed PCA analysis on the datasets and plotted the region that contains 60% of the data points to eliminate the instability caused by outliers. The highlighted points are the top 100 highest-reward points. The positions of these high-reward points varied significantly across datasets, even within the same environment and the shifting trade suggests the gradual converging process. The differences in the plot imply that for any given dataset, the data from another expertise level often lies in an OOD region. The dataset distributions of other environments are show in Appendix J. Based on this observation, we conducted OOD experiments by training a guidance network on one dataset, freezing it, and then using it to guide the training of a main diffusion model on a different dataset. Across 36 tested models, 28/36(78%) exhibited significant performance degradation. Among the eight cases with good performance, half benefited

432 from guidance modules pretrained using median replay models, which cover a broader range of data
433 and thus generalize better than other datasets.

434 From these results, we conclude that the decouple method is highly sensitive to distribution shift.
435 Training a guidance network on in-distribution data is critical for success. Otherwise, even expert-level
436 data cannot reliably guide training on a different distribution. This abolition study is consistent
437 with the conclusion in the Methodology part that OOD pretraining cannot guarantee performance
438 improvement. Another kind of limitation of Q-guided diffusion, not only in our methods but also is
439 whether the guidance should exist in the training. As shown in Table 13, for high-dimensional tasks
440 with narrow data coverage (e.g., Adroit), even slight guidance can push the policy out-of-distribution,
441 causing performance to drop significantly. In these cases, relying solely on behavioral cloning
442 proves more stable and effective. Conversely, for tasks with moderate difficulty and abundant data,
443 Q-guidance can successfully enhance performance. (details shown in the Appendix J)

444 6 RELATED WORK

445 **Modular and Decoupled Training:** Modular training has long been pursued as a desirable paradigm.
446 Recently, energy-based guidance Lu et al. (2023) was proposed, where a diffusion model is trained
447 first and an energy model is subsequently learned to provide guidance. However, such post-hoc
448 modularization has proven fragile in practice—e.g., in our own experiments, more than half of the
449 runs diverged. In contrast, our method inverts this order: we first train a guidance module using
450 supervised learning from offline data, and then use this frozen module to learn a diffusion. It serves
451 as a general-purpose enhancement to existing architectures of CFG. Another line of research is
452 semi-modularized: although IQL can be seen as a modular paradigm—learning Q-values first and
453 then selecting one-hot action in the inference stage—it does not apply guidance during training,
454 and its performance often lags behind joint methods such as DQL. Classic offline RL algorithms
455 such as IQL Kostrikov et al. (2022), CQL Kumar et al. (2020), and BRAC Wu et al. (2019) have
456 also emphasized batch-constrained or conservative regularization, demonstrating the importance
457 of stability under distributional shift. Our work proposes a guidance-first modular framework that
458 enhances training in offline RL, and it provides the first systematic examination of when modular
459 guidance is feasible and demonstrates that, under the challenges of offline RL, modular training can
460 succeed if designed around guidance-first principles. **semble-based Error Correction:** It is inspired
461 by Double Q-learning and ensemble learning, we leverage two independently initialized modules
462 to cancel stochastic biases inherent in a single model, a technique widely adopted in reinforcement
463 learning but underexplored in diffusion-based policies. **Plug-and-Play Modular Composition.** The
464 idea of plug-and-play composability is to treat pretrained modules—originally developed for other
465 purposes—as reusable building blocks. Such composability is rarely studied and highly empirical.
466 We propose that plug-and-play composition requires distributional alignment. There is an exhaustive
467 related work section in Appendix K.

468 While our current framework adopts existing diffusion formulations, we believe the paradigm of
469 pretraining guidance models holds great potential for future research. Pretrained guidance not only
470 enforces in-distribution sampling and constrains out-of-distribution behaviors but also facilitates
471 modularity and adaptability. We envision that future diffusion algorithms specifically designed to
472 accommodate pretrained guidance could further improve efficiency, reduce variance, and accelerate
473 convergence.

474 REPRODUCIBILITY AND ETHICS STATEMENT

475 Our work builds on the Clean Diffusion Wang et al. (2024) framework with modular modifications.
476 Experiments use three random seeds; hyperparameters are in Appendix F. Code, pretrained
477 models, and scripts are available at [https://github.com/juliachen2357-maker/
478 modular-diffusion-policy](https://github.com/juliachen2357-maker/modular-diffusion-policy). The study uses public benchmarks only and has no foreseeable
479 negative societal impact, though safe and fair deployment should be considered.

482 REFERENCES

483 Dongjun Kim, Yeongmin Kim, Se Jung Kwon, Wanmo Kang, and Il-Chul Moon. Refining generative
484 process with discriminator guidance in score-based diffusion models. In *International Conference*
485

- 486 *on Machine Learning*, volume 202, pages 1–25. PMLR, 2023.
- 487
- 488 Sergey Levine, Aviral Kumar, George Tucker, and Justin Fu. Offline reinforcement learning: Tutorial,
489 review, and perspectives on open problems. *arXiv preprint arXiv:2005.01643*, 2020.
- 490 Michael Janner, Yilun Du, Joshua B. Tenenbaum, and Sergey Levine. Planning with diffusion
491 for flexible behavior synthesis. In *International Conference on Machine Learning*, 2022. URL
492 <https://arxiv.org/abs/2205.09991>.
- 493
- 494 Bingyi Kang, Haotian Shi, Yilun Tan, Mingming Zhu, Dahuan Lin, and Hang Zhou. Efficient
495 diffusion policies for offline reinforcement learning. *arXiv preprint arXiv:2310.03573*, 2023.
- 496 Philippe Hansen-Estruch, Ilya Kostrikov, Michael Janner, Jakub Grudzien Kuba, and Sergey Levine.
497 Idql: Implicit q-learning as an actor-critic method with diffusion policies. *arXiv preprint*
498 *arXiv:2304.10573*, 2023.
- 499
- 500 Hado van Hasselt. Double q-learning. In *Advances in Neural Information Processing Systems*,
501 volume 23, pages 2613–2621, 2010.
- 502 Anurag Ajay, Yilun Du, Abhi Gupta, Joshua B. Tenenbaum, Tommi S. Jaakkola, and Pulkit Agrawal.
503 Is conditional generative modeling all you need for decision making? In *International Conference*
504 *on Learning Representations*, 2023. Top 5% of submissions.
- 505
- 506 Scott Fujimoto, Edoardo Conti, Mohammad Ghavamzadeh, and Joelle Pineau. Benchmarking batch
507 deep reinforcement learning algorithms. *arXiv preprint arXiv:1910.01708*, 2019a.
- 508 Takeru Miyato, Toshiki Kataoka, Masanori Koyama, and Yuichi Yoshida. Spectral normalization
509 for generative adversarial networks. In *Advances in Neural Information Processing Systems*,
510 volume 31, 2018. URL [https://proceedings.neurips.cc/paper/2018/file/](https://proceedings.neurips.cc/paper/2018/file/1802.05957-Paper.pdf)
511 [1802.05957-Paper.pdf](https://proceedings.neurips.cc/paper/2018/file/1802.05957-Paper.pdf).
- 512
- 513 Justin Fu, Aviral Kumar, Ofir Nachum, George Tucker, and Sergey Levine. D4rl:
514 Datasets for deep data-driven reinforcement learning. In *Advances in Neural*
515 *Information Processing Systems Datasets and Benchmarks Track*, 2021. URL
516 [https://datasets-benchmarks-proceedings.neurips.cc/paper/2021/](https://datasets-benchmarks-proceedings.neurips.cc/paper/2021/file/7b5b23f4aadf9513306bcd59afb6e4c9-Paper-round2.pdf)
517 [file/7b5b23f4aadf9513306bcd59afb6e4c9-Paper-round2.pdf](https://datasets-benchmarks-proceedings.neurips.cc/paper/2021/file/7b5b23f4aadf9513306bcd59afb6e4c9-Paper-round2.pdf).
- 518 The Farama Foundation. Gymnasium: A standard api for reinforcement learning environments.
519 <https://github.com/Farama-Foundation/Gymnasium>, 2024.
- 520
- 521 Zhendong Wang, Xinyang Zhang, Bowen Liu, et al. Cleandiffuser: Diffusion library and benchmark
522 for decision-making. In *Advances in Neural Information Processing Systems Datasets and*
523 *Benchmarks Track*, 2024. URL <https://arxiv.org/abs/2406.09509>.
- 524 Cheng Lu, Shengming Ruan, Yang Song, and Stefano Ermon. Energy-based guidance for diffusion
525 models. In *International Conference on Machine Learning*, 2023.
- 526
- 527 Ilya Kostrikov, Ashvin Nair, and Sergey Levine. Offline reinforcement learning with implicit
528 q-learning. In *International Conference on Learning Representations (ICLR)*, 2022.
- 529 Aviral Kumar, Aurick Zhou, George Tucker, and Sergey Levine. Conservative q-learning for offline
530 reinforcement learning. In *Advances in Neural Information Processing Systems*, volume 33, pages
531 1179–1191, 2020.
- 532
- 533 Yifan Wu, George Tucker, and Ofir Nachum. Behavior regularized offline reinforcement learning. In
534 *Advances in Neural Information Processing Systems*, pages 7968–7978, 2019.
- 535 Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. In *Advances in*
536 *Neural Information Processing Systems*, volume 33, pages 6840–6851, 2020.
- 537
- 538 Xingyu Ma, Yuchen Yang, Yifei Chen, Han Liu, and Tong Zhang. Dice: Offline reinforcement
539 learning with diffusion models. In *International Conference on Learning Representations (ICLR)*,
2024.

- 540 Q. Miao et al. Learning with noisy labels using collaborative sample selection. In *Advances in Neural*
541 *Information Processing Systems*, 2023.
- 542 Prafulla Dhariwal and Alex Nichol. Diffusion models beat gans on image synthesis. In *Advances in*
543 *Neural Information Processing Systems*, 2021.
- 544 Scott Fujimoto, David Meger, and Doina Precup. Off-policy deep reinforcement learning without
545 exploration. In *International Conference on Machine Learning*, pages 2052–2062, 2019b.
- 546 Tianhe Yu, Aviral Kumar, Yevgen Chebotar, Karol Hausman, Gaurav Sukhatme, Sergey Levine, and
547 Chelsea Finn. Mopo: Model-based offline policy optimization. In *Advances in Neural Information*
548 *Processing Systems*, pages 14129–14142, 2020.
- 549 Ashvin Nair, Aviral Kumar, Chelsea Finn, and Sergey Levine. Accelerating online reinforcement
550 learning with offline datasets. In *Advances in Neural Information Processing Systems*, 2020.
- 551 Scott Fujimoto, Herke van Hoof, and David Meger. Addressing function approximation error in
552 actor-critic methods. In *International Conference on Machine Learning*, pages 1582–1591, 2018.
- 553 Ian Osband, Charles Blundell, Alexander Pritzel, and Benjamin Van Roy. Deep exploration via
554 bootstrapped dqn. In *Advances in Neural Information Processing Systems*, volume 29, pages
555 4026–4034, 2016.
- 556 Oron Anschel, Nir Baram, and Nahum Shimkin. Averaged-dqn: Variance reduction and stabilization
557 for deep reinforcement learning. In *International Conference on Machine Learning*, pages 176–185,
558 2017.
- 559 Qingfeng Lan, Yuandong Pan, Ling Pan, Zhongxiang Zhou, Bin Hu, and Zheng Wang. Maxmin
560 q-learning: Controlling the estimation bias of q-learning. In *Advances in Neural Information*
561 *Processing Systems*, volume 33, pages 14157–14168, 2020.
- 562 Alex Nichol, Prafulla Dhariwal, Aditya Ramesh, Pranav Shyam, Pam Mishkin, Bob McGrew, Ilya
563 Sutskever, and Mark Chen. Glide: Towards photorealistic image generation and editing with
564 text-guided diffusion models. In *Advances in Neural Information Processing Systems*, 2021.
- 565 Aditya Ramesh, Mikhail Pavlov, Gabriel Goh, Scott Gray, Chelsea Voss, Alec Radford, Mark Chen,
566 and Ilya Sutskever. Zero-shot text-to-image generation. In *International Conference on Machine*
567 *Learning*, 2021.
- 568 Jacob Andreas, Dan Klein, and Sergey Levine. Modular multitask reinforcement learning with policy
569 sketches. In *International Conference on Machine Learning*, 2017. URL <https://arxiv.org/abs/1704.06643>.
- 570 Xue Bin Peng, Angjoo Kanazawa, Sam Toyer, Pieter Abbeel, and Sergey Levine. Composing
571 complex skills by learning transition policies. In *IEEE International Conference on Robotics and*
572 *Automation*, 2019. URL <https://arxiv.org/abs/1906.01068>.
- 573 Chong Mou, Yabo Zhang, Yixiao Li, et al. T2i-adapter: Learning adapters to dig out more controllable
574 ability for text-to-image diffusion models. *arXiv preprint arXiv:2302.08453*, 2023.
- 575 Lvmin Zhang, Yongjie Rao, and Maneesh Agrawala. Adding conditional control to text-to-image
576 diffusion models. In *IEEE International Conference on Computer Vision*, 2023.
- 577 Zexiang Ye, Yixiao Zhang, et al. Ip-adapter: Text-compatible image prompt adapter for text-to-image
578 diffusion models. *arXiv preprint arXiv:2308.06721*, 2023.
- 579 Patrick Lewis, Ethan Perez, Aleksandra Piktus, Fabio Petroni, Vladimir Karpukhin, Naman Goyal,
580 Heinrich Küttler, Mike Lewis, Wen-tau Yih, Tim Rocktäschel, Sebastian Riedel, and Douwe
581 Kiela. Retrieval-augmented generation for knowledge-intensive nlp tasks. In *Advances in Neural*
582 *Information Processing Systems*, 2020.
- 583 Kelvin Guu, Kenton Lee, Zora Tung, Panupong Pasupat, and Ming-Wei Chang. Realm:
584 Retrieval-augmented language model pre-training. In *International Conference on Machine*
585 *Learning*, pages 3929–3938, 2020.

594
595
596
597
598
599
600
601
602
603
604
605
606
607
608
609
610
611
612
613
614
615
616
617
618
619
620
621
622
623
624
625
626
627
628
629
630
631
632
633
634
635
636
637
638
639
640
641
642
643
644
645
646
647

Gautier Izacard and Edouard Grave. Leveraging passage retrieval with generative models for open domain question answering. In *Conference on Empirical Methods in Natural Language Processing*, pages 8741–8754, 2021.

Gautier Izacard, Edouard Grave, Fabio Petroni, Lucas Hosseini, Timo Schick, Pierre-Emmanuel Mazare, Julian Eisenschlos, Sebastian Petrov, Vladimir Karpukhin, Patrick Lewis, Wen-tau Yih, Tim Rocktaschel, Sebastian Riedel, and Douwe Kiela. Atlas: Few-shot learning with retrieval augmented language models. In *International Conference on Machine Learning*, 2022.

A REMARKS

In Table 4, we analyze the number of parameters in the guidance module and the diffusion model. The parameter counts objectively reflect the memory footprint and computational cost of each component. Since we pre-train the guidance module, the peak memory usage during GFDT training only needs to accommodate the largest of these components. Therefore, we consider the effective reduction to be determined by the smaller of the two overlapping components, which corresponds to a 38.1% decrease for our model.

B Q-GUIDED DIFFUSION

Q-Guided Diffusion Janner et al. (2022) has a forward noising process and a reverse denoising process. The forward noising process gradually perturbs a clean action a_0 into a noisy version a_k :

$$a_k = \sqrt{\bar{\alpha}_k} a_0 + \sqrt{1 - \bar{\alpha}_k} \epsilon, \quad \epsilon \sim \mathcal{N}(0, I), \quad \bar{\alpha}_k = \prod_{i=1}^k \alpha_i,$$

where $k \in \{1, \dots, T\}$ denotes the diffusion timestep (or noise level), and $\bar{\alpha}_k$ is the cumulative product of noise scheduling coefficients. A denoising network $\epsilon_\theta(a_k, k, s)$ is trained to recover the clean action a_0 from its noisy counterpart $a_k \sim \mathcal{N}(0, I)$. During reverse diffusion, we can form an estimate of the clean action:

$$\hat{a}_0 = \frac{1}{\sqrt{\bar{\alpha}_k}} \left(a_k - \sqrt{1 - \bar{\alpha}_k} \cdot \epsilon_\theta(a_k, k, s) \right). \quad (6)$$

The behavior clone loss is formulated as: $\mathcal{L}_{\text{BC}} = \mathbb{E}_{a_0, \epsilon, t} [\|\epsilon_\theta(a_t, t, s) - \epsilon\|^2]$. To incorporate reward information, Q-guided diffusion introduces an actor loss \mathcal{L}_Q , and the entire loss becomes: $\mathcal{L}_{\text{actor}} = \mathcal{L}_{\text{BC}} + \eta \mathcal{L}_Q$, which is the sum of the behavior cloning loss, and Q-loss, which is controlled by the strength of η . Here $Q_\phi(s, a)$ denotes the learned Q-function, i.e., an estimation of the expected discounted return starting from state s and action a : $Q_\phi(s, a) \approx \mathbb{E}[\sum_{t=0}^{\infty} \gamma r_t \mid s_0 = s, a_0 = a]$, (γ is the discount factor, r_t the reward at step t). For stable guidance, we normalize this value by its expectation over sampled actions, defining $\tilde{Q}_\phi(s, a) = Q_\phi(s, a) / \mathbb{E}_a[Q_\phi(s, a)]$. This normalization ensures scale invariance of the guidance signal and enables interchangeable use of guidance modules across actor-critic variants (e.g., DQL and IDQL) without a mismatch in magnitude. (A more detailed derivation is provided in section 3.)

A similar design also exists at inference time, the Q-network provides explicit guidance by adjusting the sampled actions:

$$\tilde{a}_0 = a_0 + \lambda \frac{\nabla_a Q_\phi(s, a_0)}{\|\nabla_a Q_\phi(s, a_0)\| + \epsilon}, \quad (7)$$

where λ balances exploitation of the learned Q-function against fidelity to the diffusion. The normalizing process is inherent from Ho et al. (2020) and is expected to be important in the modular and plug-and-play designs.

C ALGORITHMIC DETAILS

C.1 GFDT ALGORITHMIC DETAILS

The entire training and inference process of GFDT is shown in Fig.4, as well as the algorithm 1.

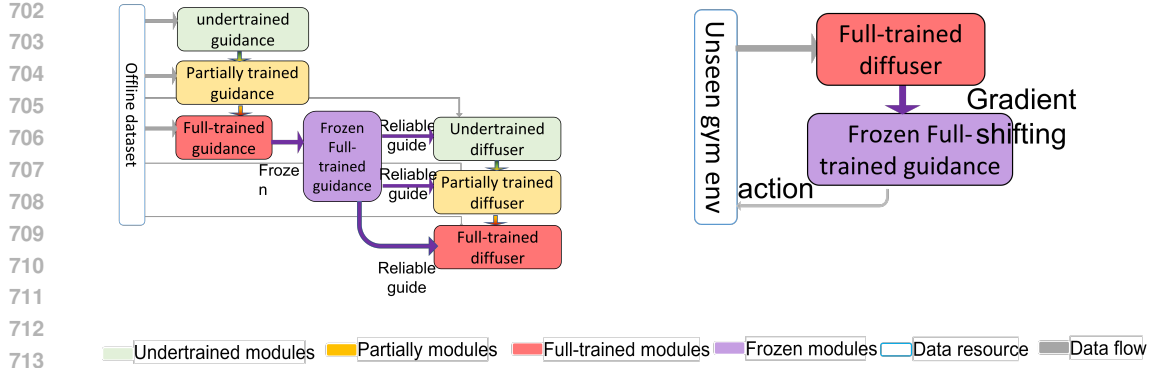


Figure 4: Training and inference stage of GFDT

Algorithm 1 Guidance-First Diffusion Training (GDFT)

Require: Offline dataset $\mathcal{D} = \{(s, a, r, s')\}$, Q-network Q_ϕ , diffusion model ϵ_θ , guidance scale λ , training steps N_q, N_θ

1: // **Step 1: Train Q-function (guidance model)**

2: **for** $i = 1$ to N_q **do**

3: Sample minibatch $(s, a, r, s') \sim \mathcal{D}$

4: TD target: $y \leftarrow r + \gamma \cdot \max_{a'} Q_\phi(s', a')$

5: Update Q_ϕ to minimize $\mathcal{L}_Q = \|Q_\phi(s, a) - y\|^2$

6: Freeze Q_ϕ

7: // **Step 2: Train diffusion model with frozen guidance**

8: **for** $j = 1$ to N_θ **do**

9: Sample $(s, a_0) \sim \mathcal{D}$

10: Add noise: $a_t \leftarrow \sqrt{\bar{\alpha}_t} a_0 + \sqrt{1 - \bar{\alpha}_t} \cdot \epsilon, \epsilon \sim \mathcal{N}(0, I)$

11: Predict noise: $\hat{\epsilon} \leftarrow \epsilon_\theta(a_t, s, t)$

12: Update ϵ_θ to minimize $\mathcal{L}_{\text{diff}} = \|\hat{\epsilon} - \epsilon\|^2 + \mathcal{L}_Q$

▷ Q_ϕ frozen; \mathcal{L}_Q still updates ϵ_θ

13: **Return** trained ϵ_θ , frozen Q_ϕ

14:

15: **Inference:**

16: **for** denoise steps **do**

17: Sample candidate $a_k \sim \pi_\theta(\cdot | s)$

18: Apply guidance: $a_{k-1} \leftarrow a_k + \lambda \cdot \nabla_a Q_\phi(s, a_k)$

Important explanation of the frozen Q network still updating the diffusion network The details of GFDT have been thoroughly explained in the main content. A reasonable concern is the role of Q_ϕ . Although Q_ϕ is frozen and does not update the parameters of the guidance network, its output values are still used to adjust the parameters of the diffusion network, encouraging it to generate samples with higher Q values.

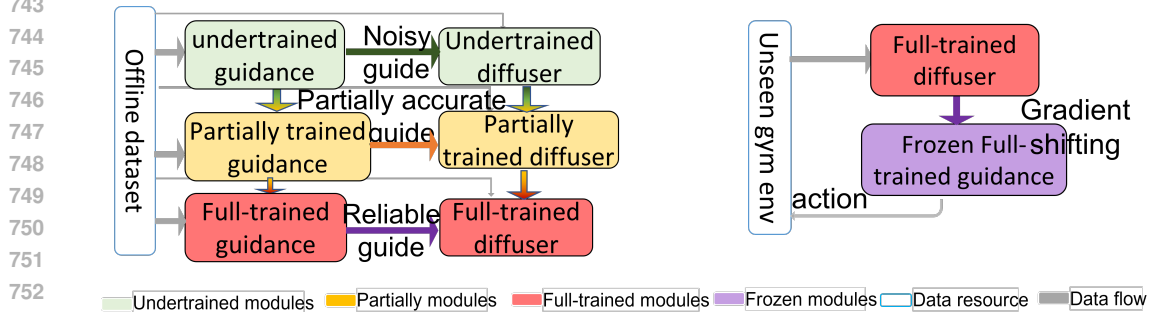


Figure 5: Training and inference stage of traditional CFG

Algorithm 2 Traditional Gradient Guidance Training

Require: Offline dataset $\mathcal{D} = \{(s, a, r, s')\}$, diffusion model ϵ_θ , Q-function Q_ϕ (learned jointly), training steps N_θ

- 1: **for** $j = 1$ to N_θ **do**
- 2: Sample $(s, a_0, r, s') \sim \mathcal{D}$
- 3: TD target: $y \leftarrow r + \gamma \cdot \max_{a'} Q_\phi(s', a')$
- 4: Update Q_ϕ to minimize $\mathcal{L}_Q = \|Q_\phi(s, a_0) - y\|^2$
- 5: Add noise: $a_t \leftarrow \sqrt{\bar{\alpha}_t} a_0 + \sqrt{1 - \bar{\alpha}_t} \cdot \epsilon, \quad \epsilon \sim \mathcal{N}(0, I)$
- 6: Predict noise: $\hat{\epsilon} \leftarrow \epsilon_\theta(a_t, s, t)$
- 7: Update ϵ_θ to minimize $\mathcal{L}_{\text{diff}} = \|\hat{\epsilon} - \epsilon\|^2 + \mathcal{L}_Q$
- 8: **Return** trained ϵ_θ , jointly-trained Q_ϕ
- 9:
- 10: **Inference:**
- 11: **for** denoise steps **do**
- 12: Sample candidate $a_k \sim \pi_\theta(\cdot|s)$
- 13: Apply guidance: $a_{k-1} \leftarrow a_k + \lambda \cdot \nabla_a Q_\phi(s, a_k)$
- 14: **Return** a^*

C.2 IDQL ALGORITHMIC DETAILS

Training. Implicit Diffusion Q-learning (IDQL) decouples the training of the diffusion model from the Q-value estimator. The diffusion policy $\pi_\theta(a|s)$ is trained purely by behavior cloning (BC) from the offline dataset $\mathcal{D} = \{(s, a)\}$, i.e.,

$$\min_{\theta} \mathbb{E}_{(s,a) \sim \mathcal{D}} [\|a - \pi_\theta(s)\|^2], \quad (8)$$

without any reward or Q-guidance incorporated into the diffusion process. In parallel, a separate Q-network $Q_\phi(s, a)$ is learned by standard temporal-difference (TD) regression:

$$\min_{\phi} \mathbb{E}_{(s,a,r,s') \sim \mathcal{D}} \left[\left(Q_\phi(s, a) - (r + \gamma \max_{a'} Q_\phi(s', a')) \right)^2 \right]. \quad (9)$$

Notably, the Q-network is updated together with the diffusion model during training, although it does not directly affect the diffusion optimization.

Inference. At deployment, the diffusion model generates n candidate actions $\{a_1, a_2, \dots, a_n\}$ sampled from the diffusion model $\pi_\theta(\cdot|s)$. These are then evaluated with the learned Q-network, and the action with the highest Q-value is selected:

$$a^* = \arg \max_{i=1, \dots, K} Q_\phi(s, a_i). \quad (10)$$

This procedure can also be interpreted as a one-hot selection mechanism over the candidate actions, where the Q-network acts as a ranking function. The algorithm of IDQL is in

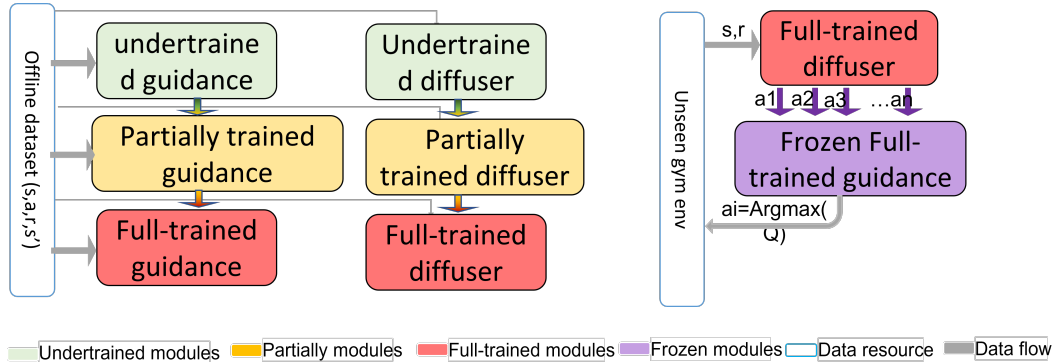


Figure 6: Training and inference stage of IDQL

Comparison between GFDT and IDQL The key distinction between IDQL and our proposed GFDT lies in how Q-guidance is incorporated. In IDQL, a separate Q-guidance network is trained, but

Algorithm 3 Implicit Diffusion Q-learning (IDQL)

```

810
811 1: Input: offline dataset  $\mathcal{D}$ 
812 2: Initialize diffusion policy  $\pi_\theta$  and Q-network  $Q_\phi$ 
813 3: while not converged do
814 4:   Sample batch  $(s, a, r, s')$  from  $\mathcal{D}$ 
815 5:   Update  $\pi_\theta$  by behavior cloning on  $(s, a)$ 
816 6:   Update  $Q_\phi$  by TD regression on  $(s, a, r, s')$ 
817
818 7: Inference:
819 8: Input: state  $s$ 
820 9: for  $k = K, K - 1, \dots, 1$  do
821 10:   for  $n = 1, \dots, N$  do
822 11:     Sample candidate action  $a_n^{(k)} \sim \pi_\theta(\cdot | s)$ 
823 12:     Evaluate  $Q_\phi(s, a_n^{(k)})$ 
824 13:     Select best action at step  $k$ :
825
826 
$$a^{(k)*} = \arg \max_{n=1, \dots, N} Q_\phi(s, a_n^{(k)})$$

827
828 14: Return  $a^{(1)*}$ 

```

reward information is not injected into the diffusion model during training. This design makes IDQL relatively stable under the distribution, since the learned policy is not directly biased by Q-values that could otherwise push the distribution toward out-of-distribution regions. However, the downside is that the policy is less reward-optimal, because it is guided primarily by behavioral cloning rather than directly exploiting reward signals. At inference time, IDQL applies reward guidance only as a one-hot selection among the generated actions. While this procedure ensures the selected actions are rewarding, they remain restricted to the support of the behaviorally cloned generator, which limits the achievable performance compared to gradient-guided methods.

In contrast, GFDT explicitly integrates Q-guidance into the generative diffusion process. Reward information actively shapes the learned distribution during training, which leads to more reward-aligned behaviors. At inference time, GFDT further applies a gradient shift, using gradient descent to nudge the generated action toward the locally optimal point. As a result, GFDT consistently outperforms both the baseline CFG and IDQL in experiments. Although gradient-based guidance has the potential to destabilize training, normalization and a staged training scheme mitigate this risk: reconstruction dominates in the early phase and Q-guidance gradually takes effect afterwards, and therefore, destabilized training can manifest a certain V-shape mode. (see Appendix J).

D EFFICIENT DIFFUSION POLICY (EDP)

Efficient Diffusion Policy (EDP) is a variant of Q-guided diffusion that primarily improves computational efficiency. Its central idea is **one-step denoising**, which bypasses the expensive multi-step reverse process used in standard diffusion policies.

Standard Diffusion vs. One-step Denoising In traditional diffusion policies, generating a clean action a^0 requires running a long reverse chain:

$$a^T \rightarrow a^{T-1} \rightarrow \dots \rightarrow a^0,$$

where T is typically large (e.g., 100–1000). This iterative procedure is computationally expensive.

Instead, EDP corrupts a dataset action a^0 directly into a noisy action a^k in one step:

$$a^k = \sqrt{\bar{\alpha}_k} a^0 + \sqrt{1 - \bar{\alpha}_k} \epsilon, \quad \epsilon \sim \mathcal{N}(0, I). \quad (11)$$

Then, it reconstructs an approximate clean action \hat{a}^0 by applying the denoiser once:

$$\hat{a}^0 = \frac{1}{\sqrt{\bar{\alpha}_k}} a^k - \frac{\sqrt{1 - \bar{\alpha}_k}}{\sqrt{\bar{\alpha}_k}} \cdot \epsilon_\theta(a^k, k, s), \quad (12)$$

where ϵ_θ is the learned denoising network conditioned on state s .

Role of Q-guidance Since the one-step approximation introduces bias, EDP relies more strongly on the Q-function to refine \hat{a}^0 . In practice, the algorithm evaluates candidate actions using Q-values and adjusts them toward high-value regions, which compensates for the reduced accuracy of the denoising step.

The main differences are 1)Cost reduction: iterative sampling ($O(T)$ steps) \rightarrow one-step approximation ($O(1)$). 2)Trade-off:less precise denoising \rightarrow stronger dependence on Q-guidance. 3)Outcome: training time reduced by orders of magnitude (days \rightarrow hours), at the expense of slightly noisier reconstructions. In summary, EDP is an acceleration of traditional CFG diffusion.

E TRAJECTORY-BASED GUIDANCE: WHY IT FAILS IN OFFLINE RL

This algorithm currently does not apply to trajectory-based methods. We trained several trajectory-return guided variants; however, the guidance estimates did not *converge* in the actionable sense: they moved from near-zero to a coarse range (e.g., 200–300 in a median expert dataset with Q values of $\approx 200-300$) but could not refine within that band, making the guidance ineffective. If the guidance itself is not converged, our method also does not make a difference. We attribute this to (i) noisy returns and (ii) long-horizon credit assignment, both exacerbated in offline settings without on-policy rollouts. By contrast, TD-based (s, a) -level guidance provides stable, local gradients that are compatible with diffusion updates.

F EXPERIMENTAL SETUP AND HYPERPARAMETERS

.All models are trained using the D4RLMuJoCoTD Dataset Fu et al. (2021). It provides pre-collected trajectories of varying quality, including expert, medium-expert, medium-replay, and medium datasets, enabling rigorous training of offline RL algorithms under diverse data distributions. The evaluation is done in randomly initialized environments. All the gradient steps mentioned are with respect to a batch size of 256.

Wherever possible, we adopt the original hyperparameter settings from the paper of Wang et al. (2024). Intentionally, we do not modify any training-related components—including the optimizer, learning rate, batch size, architecture, or loss function. Because a key strength of our method is that it achieves superior performance without requiring any changes to the other parameters other than pretraining or modularize. This highlights the robustness and plug-and-play nature of our approach. Final results are reported, averaged over 50 evaluation episodes. Performance is measured by normalized return, and we report both the mean(Section 4.1) and variance deduction(Section 4.3).

All experiments were implemented in PyTorch and based on the CLEANDIFFUSER framework Wang et al. (2024). We strictly followed the hyperparameter settings from the baseline implementations to ensure fair comparison. Table 7 summarizes the key values.

Table 7: Hyperparameters used in our experiments (inherited from CleanDiffuser).

Parameter	Value	Notes
Optimizer	Adam	
Learning Rate	3×10^{-4}	fixed across all models
Batch Size	256	
Discount Factor γ	0.99	
Noise Schedule	cosine	unless otherwise specified
Number of Diffusion Steps T	5,10,20,30,40	
Actor Loss Weight η	1.0	scales \mathcal{L}_Q

For reproducibility, we will release full training scripts and environment configurations in our code repository.

G WHY THE GENERATED ACTION WITH THE GFDT AND THE MODULAR METHODS ARE IN DISTRIBUTION

Addressing a Key Concern One might naturally worry that even if the guidance module (e.g., the Q-function Q_ϕ) and the diffusion model are both trained solely on the offline dataset \mathcal{B} , their

combination during sampling could still produce out-of-distribution actions. The value guidance term $\nabla_a Q_\phi(a)$ may have a large magnitude or steep gradients, which could push the sampled action a_t away from the data manifold if not properly controlled. Therefore, the perturbation magnitude and the guidance magnitude are tightly controlled, and the guidance gradient is normalized :

$$\|a_{t+1} - a_t\| \leq \lambda + \mathcal{O}(\sqrt{\tau_t}), \lambda > 0 \text{ and } \lambda \rightarrow 0$$

ensuring that each sampling step stays within a small neighborhood of the current point. Since the diffusion model is trained on the dataset \mathcal{B} and generates samples close to it, and since the guidance is applied as a *soft* correction, the overall sampling trajectory remains near the support of \mathcal{B} . Thus, even when guided by Q_ϕ , the diffusion process remains effectively batch-constrained.

H COMPARISON TABLE OF GFDT

DQL baseline	GFDT (%)	D_DFGT (%)	D_BL (%)	GAI (%)	BC (%)	Unfreeze
HCEX	88 90/ 102.5%	91/ 102.6%	89/ 100.4%	68/ 77.7%	86/ 97.3%	90/ 101.5%
HCME	88 90/ 102.2%	90/ 101.7%	85/ 96.0%	90/ 101.6%	86/ 96.9%	89/ 100.6%
HCMR	67 68/ 100.8%	68/ 101.3%	68/ 101.4%	68/ 100.6%	56/ 82.9%	67/ 100.1%
HCMV	60 67/ 111.9%	69/ 115.2%	53/ 89.3%	53/ 89.1%	55/ 91.1%	61/ 101.5%
HOEX	167 165/ 99.0%	173/ 103.6%	164/ 98.2%	163/ 97.6%	165/99.0%	163/ 97.8%
HOME	168 172/102.6%	167/ 99.4%	161/ 95.8%	166/ 99.0%	167/99.3%	159/ 94.7%
HOMR	152 153/101.0%	153/ 100.9%	156/103.1%	122/ 80.5%	113/74.7%	150/ 99.3%
HOMV	144 147/102.2%	145/ 101.0%	133/ 92.1%	71/ 49.4%	118/82.0%	145/100.5%
WAEX	117 120/102.6%	120/ 102.1%	115/ 98.5%	120/102.3%	108/92.2%	118/100.6%
WAME	118 118/ 100%	113/ 95.9%	114/ 96.5%	115/ 97.8%	117/99.8%	119/100.7%
WAMR	93 94/ 101.4%	93/ 100.3%	83/ 89.8%	80/ 86.4%	87/ 93.7%	96/ 102.9%
WAMV	88 88/ 100.3%	88/ 100.3%	77/ 87.5%	77/ 88.2%	82/ 93.9%	88/ 100.3%

Table 8: Comparison of baseline and different methods across environments (rounded to 2 decimals and percentages to 1 decimal).

EDP baseline	GFDT	D_BL(%)	D_GFDT(%)	GAI(%)	BC(%)	Unfreeze
HCEX	0.87 0.87/100.3%	0.86/ 99.2%	0.86/ 99.9%	0.68/ 79.1%	0.86/ 99.2%	0.86/ 99.3%
HCME	0.87 0.87/100.5%	0.87/100.2%	0.87/ 100.2%	0.90/103.4%	0.86/ 98.6%	0.87/100.1%
HCMR	0.66 0.52/ 79.6%	0.65/ 98.7%	0.64/ 97.8%	0.68/103.0%	0.56/ 84.9%	0.64/ 97.9%
HCM	0.55 0.59/108.6%	0.56/101.8%	0.55/101.2%	0.53/ 97.8%	0.55/100.0%	0.55/101.4%
HOEX	1.61 1.64/101.5%	1.62/100.2%	1.62/ 100.5%	1.63/100.9%	1.65/102.4%	1.62/100.6%
HOME	1.61 1.64/101.6%	1.61/ 99.6%	1.62/ 100.2%	1.66/103.1%	1.67/103.4%	1.64/101.4%
HOMR	1.29 1.46/113.1%	1.44/110.9%	1.40/ 107.9%	1.22/ 94.2%	1.13/ 87.5%	1.20/ 92.6%
HOM	1.41 1.42/100.3%	1.37/ 97.1%	1.39/ 98.7%	0.71/ 50.4%	1.18/ 83.6%	1.42/100.3%
WAEX	1.16 1.16/100.2%	1.16/100.1%	1.16/ 100.1%	1.20/103.3%	1.08/ 93.0%	1.16/100.2%
WAME	1.16 1.17/100.2%	1.17/100.5%	1.17/ 100.5%	1.15/ 98.9%	1.17/101.0%	1.17/100.7%
WAMR	0.85 0.84/ 98.5%	0.83/ 98.1%	0.82/ 96.7%	0.80/ 94.5%	0.87/102.5%	0.85/100.2%
WAM	0.84 0.84/100.0%	0.84/ 99.1%	0.83/ 99.1%	0.77/ 91.8%	0.82/ 97.7%	0.84/100.2%

Table 9: Comparison of baseline and different methods across environments (abbreviated, compact format).

Antmaze	baseline	GFDT	D_GFDT	D_BL	GAI
large-diverse-v2	0.63	0.91 (143.92%)	0.85 (133.68%)	0.85 (133.68%)	0.87 (136.84%)
large-play-v2	0.90	0.88 (139.68%)	0.89 (99.26%)	0.88 (97.78%)	0.89 (98.52%)
medium-diverse-v2	0.93	0.97 (154.50%)	0.89 (95.71%)	0.87 (93.57%)	0.94 (100.71%)
medium-play-v2	0.67	0.91 (144.97%)	0.92 (136.63%)	0.76 (112.87%)	0.89 (131.68%)

Table 10: Antmaze results DQL.

Antmaze	baseline	GFDT	D_GFDT	D_BL	GAI
large-diverse	0.31	0.31 (102.17%)	0.35 (113.04%)	0.53 (173.90%)	0.29 (93.48%)
large-play	0.21	0.23 (106.25%)	0.23 (106.25%)	0.27 (125.00%)	0.19 (87.50%)
medium-diverse	0.67	0.52 (77.23%)	0.43 (63.37%)	0.19 (28.70%)	0.02 (2.97%)
medium-play	0.73	1.18 (160.98%)	1.15 (156.43%)	1.10 (150.10%)	0.90 (122.78%)

Table 11: Antmaze results with EDP

Table 12: Comparison with standard offline RL baselines on D4RL tasks. Results ReBR, short for REBRAC from DICEMa et al. (2024), CQL are taken from Kumar et al. (2020), IQL from Kostrikov et al. (2022), EDP-GFDT/EDP-Double-GFDT and DQL-GFDT/Double-GFDT from our modifications of EDP and DQL.

Env	ReBR	DICE	CQL	IQL	DQL_GF	DQL_D_DF	EDP_GF	EDF_D_GF
HCME	96.6	97.3±0.6	91.6±2.8	86.7±5.3	90.2±0.4	89.8±1.0	87.2±0.0	86.9±0.0
HCMR	46.6	49.2±0.9	45.5±0.5	44.2±1.2	67.9±0.3	68.2±0.4	64.4±0.0	64.3±0.0
HCMV	47.7	60.0±0.6	44.0±5.4	47.4±0.2	67.0±0.6	69.0±0.5	59.3±0.4	55.2±0.1
HOME	112.2	112.2±0.3	105.4±6.8	91.5±14.3	172.3±1.1	167.0±0.9	163.9±0.0	161.6±0.0
HOMR	101.7	102.3±2.1	95.0±6.4	94.7±8.6	119.9±11.6	153.0±0.8	146.4±1.7	139.7±6.7
HOMV	99.0	100.2±3.2	58.5±2.1	66.2±5.7	147.1±1.1	145.4±1.4	141.6±0.0	139.3±0.1
WAME	110.8	114.1±0.5	108.8±0.7	109.6±11.0	117.6±0.3	112.9±22.4	116.6±0.0	116.9±0.0
WAMR	84.1	90.8±2.6	77.2±5.5	73.8±1.8	95.6±4.8	93.3±0.6	83.8±0.4	82.3±1.0
WAMV	85.8	89.3±1.3	72.5±0.8	78.3±8.7	87.9±0.3	87.9±0.3	84.2±0.0	83.4±0.0
L-div	70.0	91.3±3.1	14.9	47.5±9.5	90.7±5.4	84.7±6.0	31.3±6.9	34.7±6.9
L-play	74.0	85.7±4.8	15.8	39.6±5.8	88.0±5.7	89.3±5.6	22.7±6.7	22.7±6.6
M-div	90.0	68.6±8.6	53.7	70.0±10.9	97.3±4.0	89.3±5.6	52.0±7.7	42.7±6.9
M-play	87.0	72.0±6.5	61.2	71.2±7.3	91.3±5.3	92.0±5.2	118.0±10.4	114.7±9.7

I DETAILS OF APPLICATION SCOPE

This plot is a structure diagram of Double Guidance.

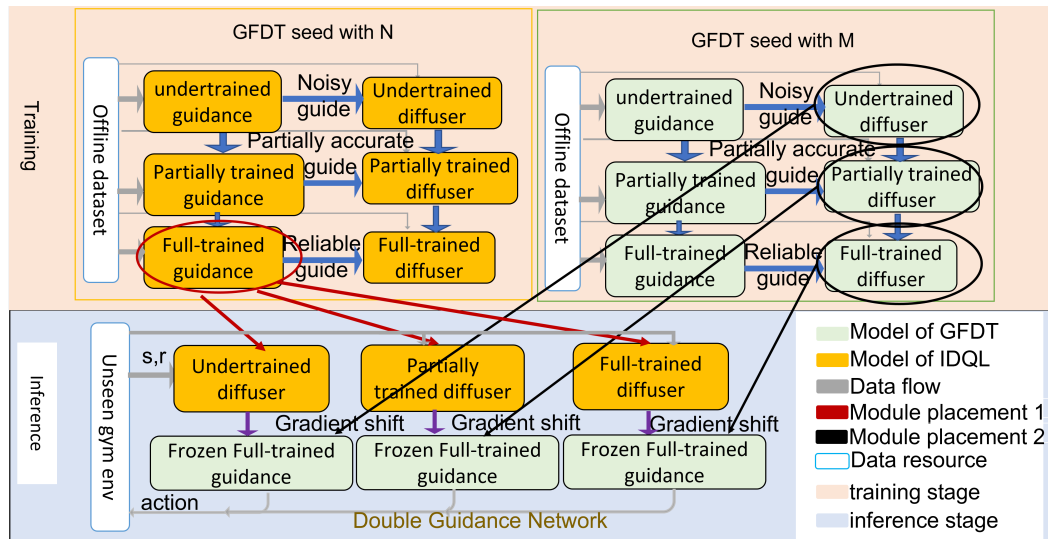
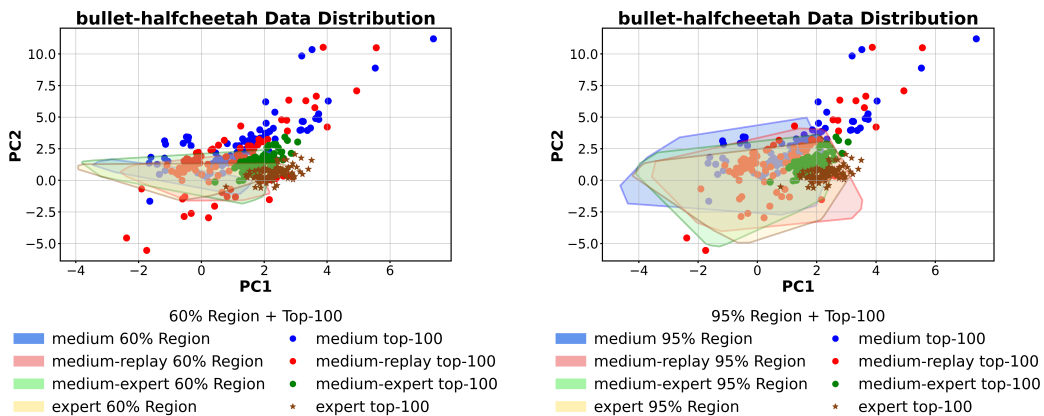


Figure 7: density hulls (top 100, $k = 20$) of Half Cheetah

J DETAILS OF APPLICATION SCOPE

1026
1027
1028
1029
1030
1031
1032
1033
1034
1035
1036
1037
1038
1039
1040
1041
1042
1043
1044
1045
1046
1047
1048
1049
1050
1051
1052
1053
1054
1055
1056
1057
1058
1059
1060
1061
1062
1063
1064
1065
1066
1067
1068
1069
1070
1071
1072
1073
1074
1075
1076
1077
1078
1079

These are all the plots of cross experiments. We did not perform data mixing experiments because, as shown in the study by Miao et al. (2023), mixing datasets generated from different policies can lead to degraded performance. Our ablation on the OOD dataset does not involve dataset mixing, since the diffusion model and the guidance model are each trained on only one dataset. Consequently, the diffusion model does not learn from multiple policies. However, when the guidance model is trained on a different dataset, it may fail to provide correct gradient signals.

(a) density hulls (top 100, $k = 20$) of Half Cheetah

(b) PCA hull and top 100 samples of Half Cheetah

1080
1081
1082
1083
1084
1085
1086
1087
1088
1089
1090
1091
1092
1093
1094
1095
1096
1097
1098
1099
1100
1101
1102
1103
1104
1105
1106
1107
1108
1109
1110
1111
1112
1113
1114
1115
1116
1117
1118
1119
1120
1121
1122
1123
1124
1125
1126
1127
1128
1129
1130
1131
1132
1133

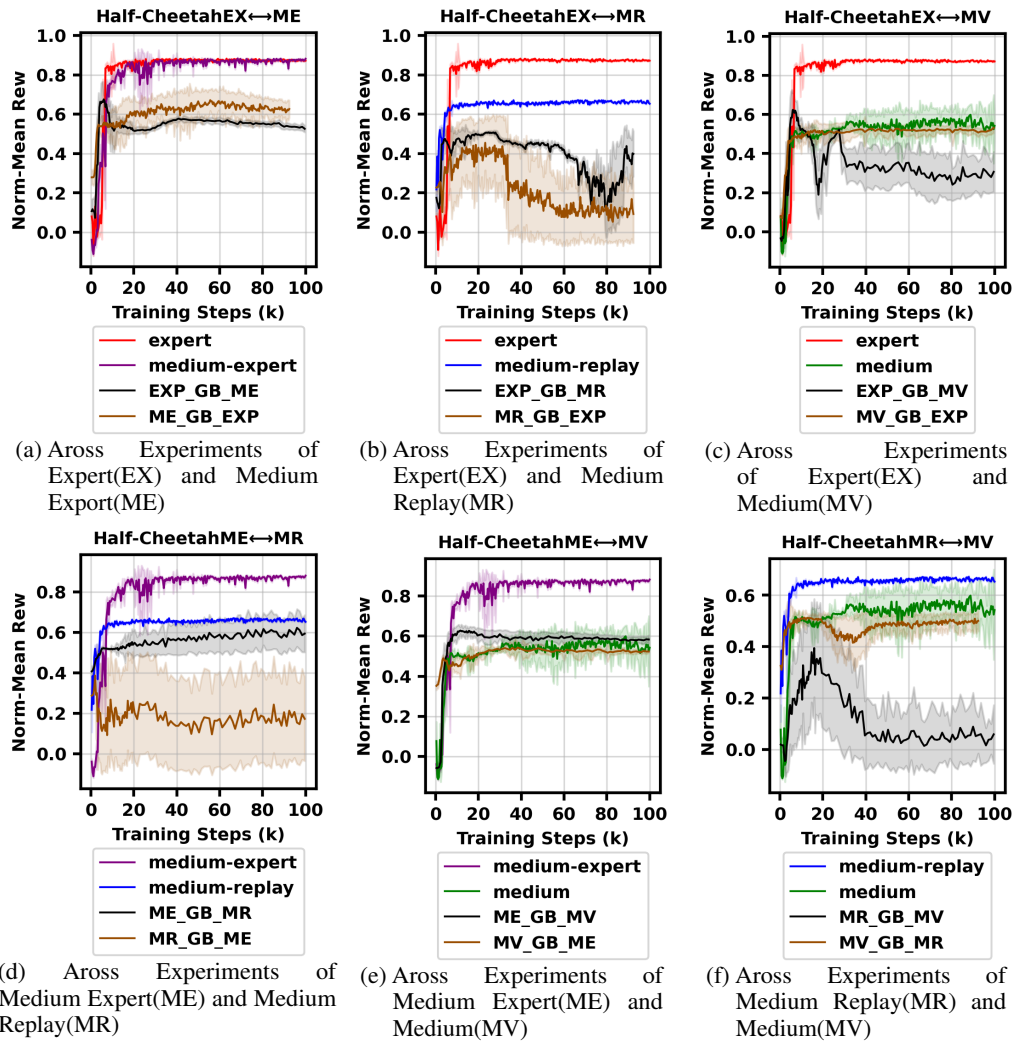
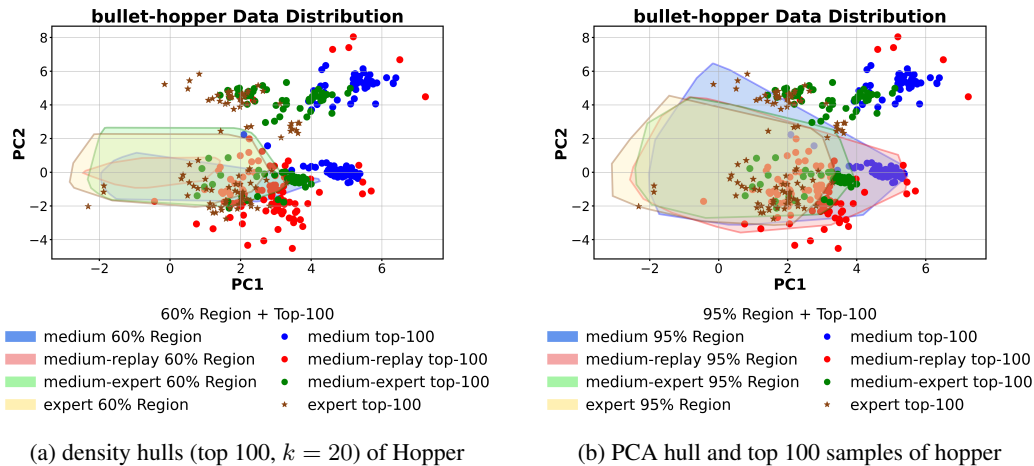


Figure 9: Half-cheetah Dataset Experiments



1134
1135
1136
1137
1138
1139
1140
1141
1142
1143
1144
1145
1146
1147
1148
1149
1150
1151
1152
1153
1154
1155
1156
1157
1158
1159
1160
1161
1162
1163
1164
1165
1166
1167
1168
1169
1170
1171
1172
1173
1174
1175
1176
1177
1178
1179
1180
1181
1182
1183
1184
1185
1186
1187

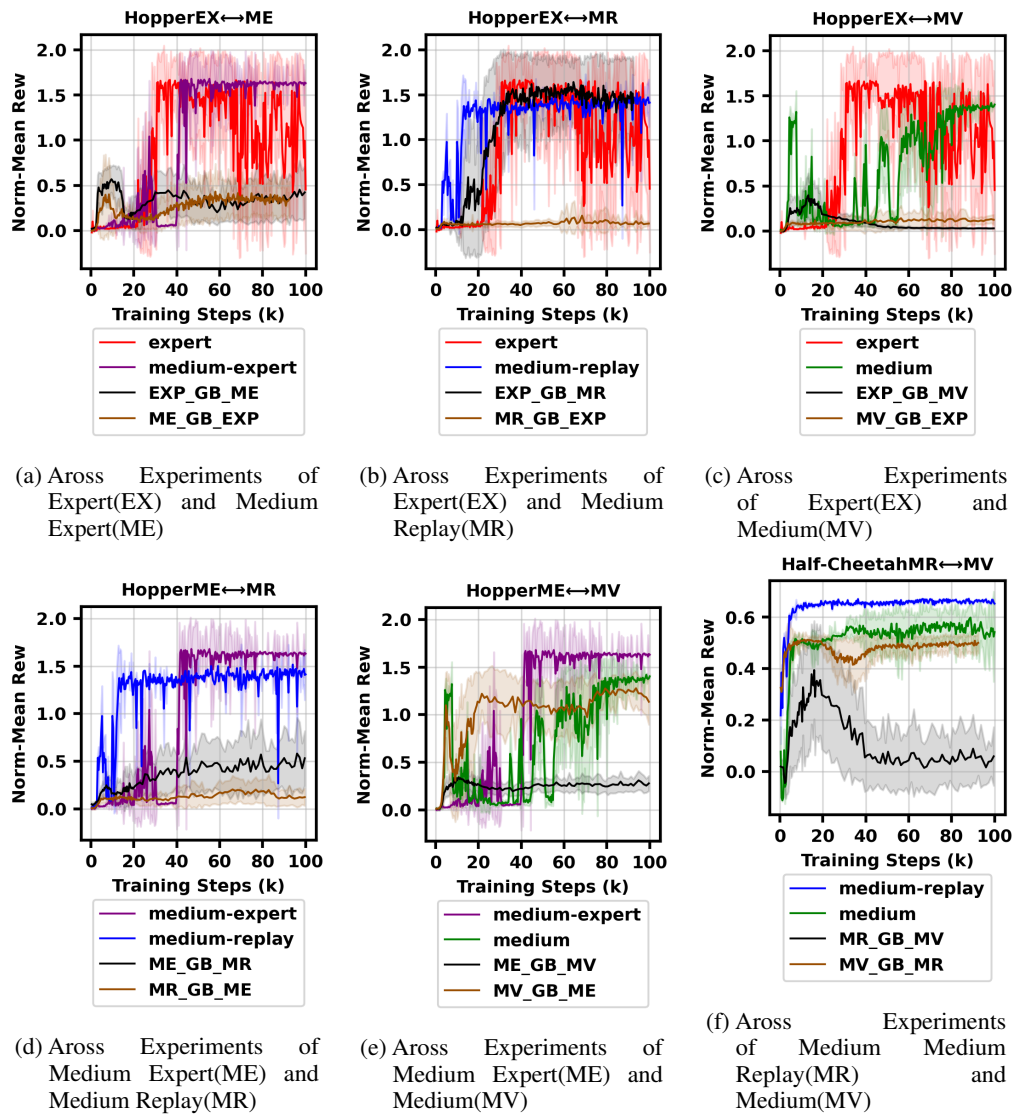
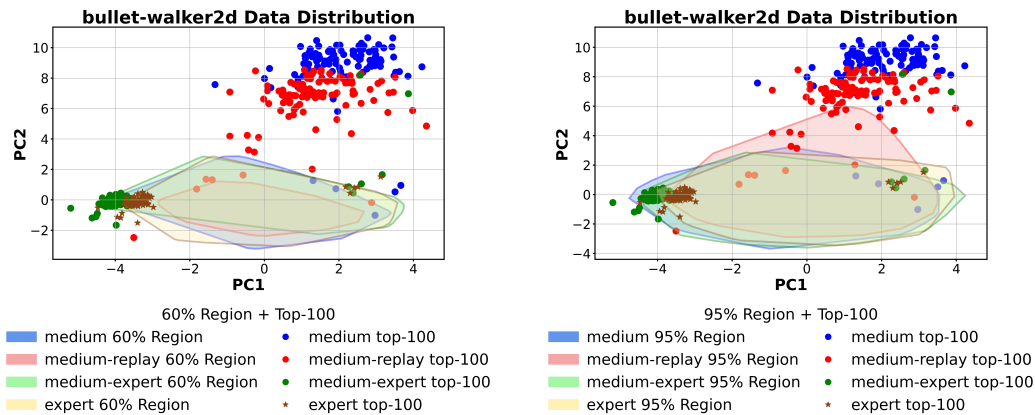


Figure 11: Hopper Dataset Experiments



(a) density hulls (top 100, $k = 20$) of Walker2d

(b) PCA hull and top 100 samples of Walker2d

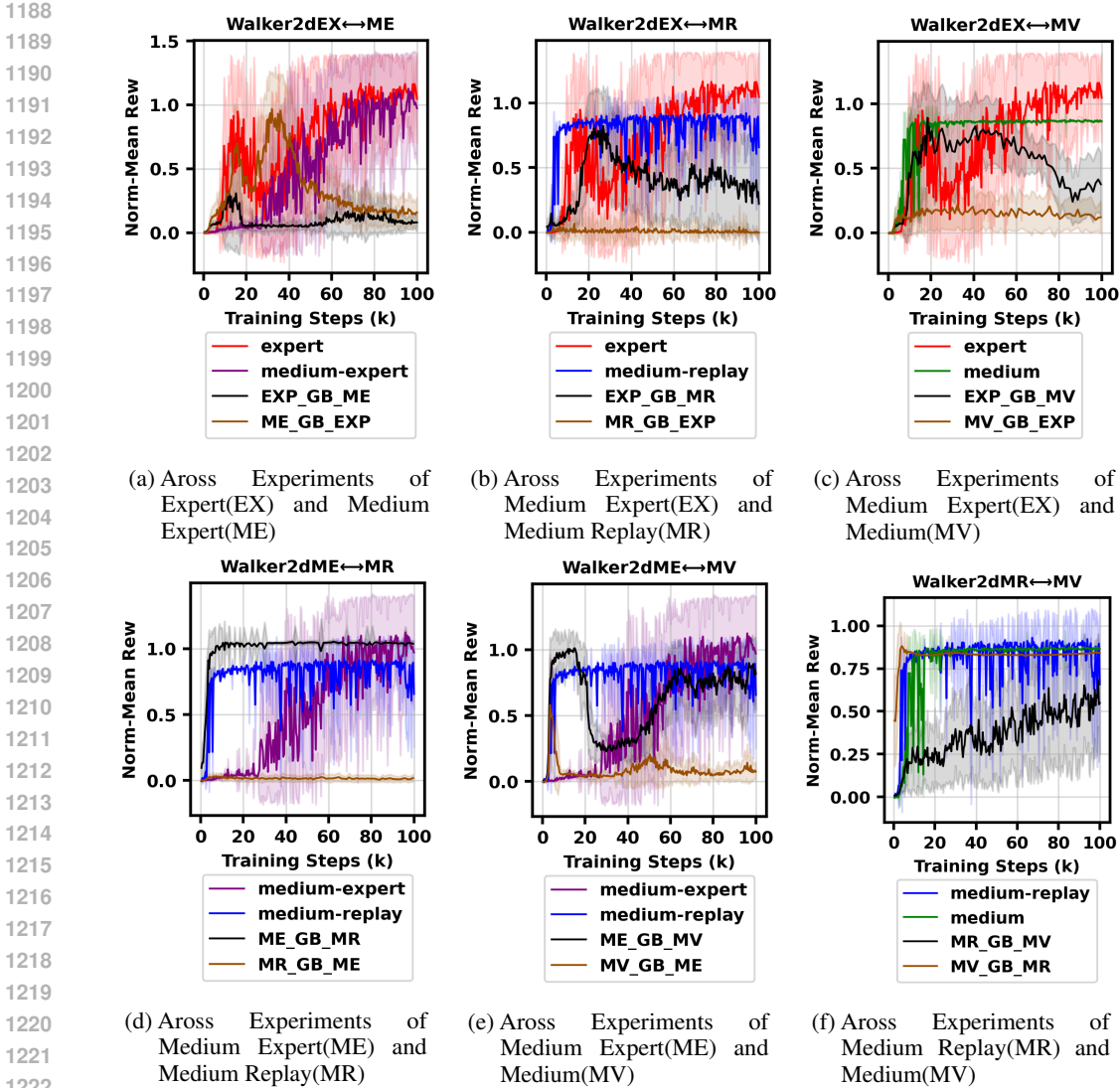


Table 13: Adroit performance with different Q-Guidance strength

Dataset	$\eta=0(BC)$	$\eta=0.1$	$\eta=1$
relocate-expert	105.99±3.78	98.20±4.61	-0.27±0.20
pen-expert	137.52±5.59	93.88±8.28	36.67±7.83
hammer-expert	124.02±3.89	119.68±4.90	0.15±2.24
door-expert-v1	105.09±9.06	102.13±3.64	-0.15±1.41

K EXTENDED DETAILS OF RELATED WORK

Modular and Decoupled Training: Modular training has long been pursued as a desirable paradigm. In the context of diffusion models, the earliest attempts at modularity can be traced to classifier guidance Dhariwal and Nichol (2021), where a pretrained diffusion model was paired with a separately trained classifier to steer sampling. This approach was unstably designed to amplify the possibility of one class while suppressing all others, ignoring that features are often shared across categories, leading to distorted and fragile guidance. Recently, energy-based guidance Lu et al. (2023) was proposed, where a diffusion model is trained first and an energy model is subsequently learned to provide guidance. However, such post-hoc modularization has proven fragile in practice—e.g., in our own experiments, more than half of the runs diverged. In contrast, our method inverts this order: we first train a guidance module using supervised learning from offline data, and then use

1242 this frozen module to learn a diffusion. It serves as a general-purpose enhancement to existing
1243 architectures of CFG. Another line of research is semi-modularized: although IQL can be seen
1244 as a modular paradigm—learning Q-values first and then selecting one-hot action in the inference
1245 stage—it does not apply guidance during training, and its performance often lags behind joint methods
1246 such as DQL. Classic offline RL algorithms such as BCQ Fujimoto et al. (2019b), CQL Kumar
1247 et al. (2020), and BRAC Wu et al. (2019) have also emphasized batch-constrained or conservative
1248 regularization, demonstrating the importance of stability under distributional shift. By contrast, our
1249 work proposes a guidance-first modular framework that enhances training in offline RL. Finally,
1250 modular training on vision or text diffusion-based generation widely exists; large language models
1251 and vision language models have Retrieval Augmented Generations or other special models separately
1252 trained as building blocks; however, since these modules use the web-harvested data, the limitations
1253 of OOD are not well-discussed. Offline RL must contend with severe out-of-distribution issues,
1254 making modularization brittle. This challenge has been extensively studied, with methods such
1255 as BRAC Wu et al. (2019), CQL Kumar et al. (2020), MOPO Yu et al. (2020), and AWAC Nair
1256 et al. (2020) proposing different strategies to mitigate extrapolation error. Our study provides the
1257 first systematic examination of when modular guidance is feasible and demonstrates that, under the
1258 challenges of offline RL, modular training can succeed if designed around guidance-first principles.

1258 **semble-based Error Correction:** Value-based reinforcement learning methods often suffer from
1259 intrinsic bias, since a single network trained on limited data can easily over- or under-estimate values
1260 in unseen regions. A classical remedy is *Double Q-learning* ??, which decouples action selection
1261 and evaluation using two networks. Similarly, *TD3* Fujimoto et al. (2018) employs twin critics and a
1262 clipped minimum operator to suppress overestimation. These methods embody the principle of error
1263 correction via redundancy, which is also central to ensemble approaches such as Bootstrapped DQN
1264 Osband et al. (2016) and Averaged-DQN Anschel et al. (2017). More recent works, e.g., *Maxmin*
1265 *Q-learning* Lan et al. (2020) and *CQL* ?, further extend this line by either using multiple critics or
1266 adding explicit penalties. Our work follows this family of ideas: inspired by Double Q-learning and
1267 ensemble learning, we leverage two independently initialized modules to cancel stochastic biases
1268 inherent in a single model, a technique widely adopted in reinforcement learning but underexplored
1269 in diffusion-based policies.

1270 **Plug-and-Play Modular Composition.** The idea of plug-and-play composability is to treat
1271 pretrained modules—originally developed for other purposes—as reusable building blocks. Such
1272 composability is rarely studied and highly empirical. We propose that plug-and-play composition
1273 requires distributional alignment. For instance, in diffusion models, CLIP-based guidance Nichol et al.
1274 (2021); Ramesh et al. (2021) applied a pretrained vision–language model as guidance, but fails on
1275 noisy intermediate states, where distribution mismatch undermines compositionality. Circumventions
1276 exist: some prior works explore modular policy composition Andreas et al. (2017); Peng et al. (2019),
1277 focusing on skill chaining or subpolicy selection. These methods apply the plug-ins as subproblem
1278 solvers in heuristic models instead of direct reward guidance. A common challenge concerns **I/O**
1279 **calibration** between modules’ signal magnitude. This issue is acute in systems with complex plug-in
1280 connections, such as adapter-based methods Mou et al. (2023); Zhang et al. (2023); Ye et al. (2023),
1281 which rely heavily on the backbone’s feature space. When transferred across architectures (e.g.,
1282 from SD1.5 to SDXL), their performance collapses due to representational mismatch, showing these
1283 adapters are not universal interfaces. The same logic applies in NLP, where Retrieval-Augmented
1284 Generation (RAG) Lewis et al. (2020) often suffers when the retriever is misaligned. Related
1285 retrieval-augmented frameworks such as REALM Guu et al. (2020), FiD Izacard and Grave (2021),
1286 and Atlas Izacard et al. (2022) further highlight the need for careful design of retriever–generator
1287 alignment and stability.

1288
1289
1290
1291
1292
1293
1294
1295

L NOTATION AND TERMINOLOGY

SYMBOLS AND VARIABLES

Symbol	Type	Meaning
\mathcal{D}	set	Offline dataset of transitions or (s, a_0) pairs
s, a, r, s'	var	State, action, reward, and next state
$\pi_\theta(a s)$	policy	Diffusion policy with parameters θ
$\epsilon_\theta(s, x_t, t)$	net	Noise predictor (denoiser)
$Q_\phi(s, a)$	net	Q-function with parameters ϕ
a_0	action	Ground-truth clean action at $t = 0$
\hat{a}_0	action	Reconstructed estimate of a_0 via ϵ_θ
x_t	action	Noisy action at diffusion step t
t, T	index	Diffusion timestep and total steps
$\alpha_t, \bar{\alpha}_t$	sched	Per-step and cumulative noise schedule coefficients
ϵ	noise	Standard Gaussian noise
η_g	scalar	Step size for one-step Q-guidance (EDP)
λ	scalar	Guidance scale during inference
L_{BC}	loss	Behavior cloning loss, $\ \epsilon - \hat{\epsilon}\ ^2$
L_Q	loss	Q-guided loss term
L_{actor}	loss	Actor loss, e.g. $L_{BC} + \eta L_Q$
γ	scalar	Discount factor in TD targets
K	count	Number of candidate actions sampled (IDQL inference)

L.1 ABBREVIATION APPENDIX

This appendix provides a comprehensive list of abbreviations used throughout the paper, particularly for environment names and model variants introduced in the experimental sections, which are not covered in the main terminology table.

D4RL BENCHMARK ENVIRONMENT ABBREVIATIONS

MODEL VARIANT ABBREVIATIONS

Abbreviation	Explanation
GFDT	Guidance-First Diffusion Training. A training paradigm where the guidance (Q-network) is pretrained and frozen before training the diffusion policy.
Baseline	The original implementation of the corresponding diffusion-based offline RL algorithm (e.g., DQL, IDQL, or EDP) without our proposed modifications.
D-Baseline	Double-Guidance Baseline. A variant of the baseline model where the guidance module used during inference is replaced with an independently initialized version of the same architecture (different random seed).
Double_GFDT	GFDT equipped with <i>Double Guidance</i> at inference time, reducing variance through independent initialization.
GAI	Guidance at Inference. A behavior cloning model trained without reward guidance during training, but augmented with Q-guidance only at inference.
Unfreeze	A variant of GFDT where the pretrained guidance module is not frozen but continues to be updated during policy training.
BC	Behavior Cloning. A supervised learning baseline that trains a policy to imitate dataset actions without reward guidance.

1350
1351
1352
1353
1354
1355
1356
1357
1358
1359
1360
1361
1362
1363
1364
1365
1366
1367
1368
1369
1370
1371
1372
1373
1374
1375
1376
1377
1378
1379
1380
1381
1382
1383
1384
1385
1386
1387
1388
1389
1390
1391
1392
1393
1394
1395
1396
1397
1398
1399
1400
1401
1402
1403

ENVIRONMENT ABBREVIATIONS

Abbreviation	Explanation
HCEX	HalfCheetah-Expert. Dataset generated by an expert policy in the HalfCheetah environment.
HCME	HalfCheetah-Medium-Expert. Dataset generated by a mixture of medium and expert policies in HalfCheetah.
HCMR	HalfCheetah-Medium-Replay. Dataset generated by replay buffer data collected from medium-performance policies.
HCMV	HalfCheetah-Medium. Dataset generated by a medium-performance policy in HalfCheetah.
HOEX	Hopper-Expert. Dataset generated by an expert policy in the Hopper environment.
HOME	Hopper-Medium-Expert. Dataset generated by a mixture of medium and expert policies in Hopper.
HOMR	Hopper-Medium-Replay. Replay buffer dataset in Hopper.
HOMV	Hopper-Medium. Dataset generated by a medium policy in Hopper.
WAEX	Walker2d-Expert. Expert policy dataset in Walker2d.
WAME	Walker2d-Medium-Expert. Mixed dataset in Walker2d.
WAMR	Walker2d-Medium-Replay. Replay dataset in Walker2d.
WAMV	Walker2d-Medium. Medium policy dataset in Walker2d.