# QUAD: Q-Gradient Uncertainty-Aware Guidance for Diffusion policies in Offline Reinforcement Learning

**Anonymous authors**Paper under double-blind review

# **ABSTRACT**

Diffusion-based offline reinforcement learning (RL) leverages Q-gradients of noisy actions to guide the denoising process. Existing approaches fall into two categories: (i) backpropagating the Q-gradient of the final denoised action through all steps, or (ii) directly estimating the Q-gradient of noisy actions. The former suffers from exploding or vanishing gradients as the number of denoising steps increases, while the latter becomes inaccurate when noisy actions deviate substantially from the dataset. In this work, we focus on addressing the limitations of the second category. We introduce QUAD, an uncertainty-aware Q-gradient guidance method. QUAD employs a Q-ensemble to estimate the uncertainty of Q-gradients and uses this uncertainty to constrain unreliable guidance during denoising. By down-weighting unreliable gradients, QUAD reduces the risk of producing suboptimal actions. Experiments on the D4RL benchmark show that QUAD outperforms state-of-the-art methods across most tasks.

#### 1 Introduction

Reinforcement learning (RL) has achieved remarkable progress in sequential decision-making tasks, ranging from games (Mnih et al., 2013; Lample & Chaplot, 2017) to robotics (He et al., 2024b; Ze et al., 2025). However, the majority of these successes rely heavily on abundant online interactions. In many real-world domains, such as healthcare, autonomous driving, and industrial control, exploration is either prohibitively costly or inherently unsafe. Offline RL addresses this challenge by learning policies purely from pre-collected datasets (Fujimoto & Gu, 2021; Zhou et al., 2025), thereby eliminating the need for online exploration. However, it suffers from distribution shift (Levine et al., 2020): the learned policy may produce actions that deviate substantially from those observed in the dataset, resulting in unreliable value estimates and degraded performance. A key contributor to this issue is the limited expressiveness of conventional policy classes (e.g. Gaussian), which struggle to capture complex, multimodal action distributions in real-world datasets, worsening the mismatch between learned and behavior policies.

Diffusion models (Ho et al., 2020) have emerged as a powerful class of policies (Chi et al., 2023), capable of capturing highly complex action distributions and generating diverse actions. Diffusion-based offline RL methods typically combines two forms of guidance: behavior cloning (BC) guidance and Q-guidance (Wang et al., 2022). BC guidance steers the denoising trajectory towards dataset-like actions, thereby alleviating distributional shift, whereas Q-guidance leverages value estimates to promote higher-quality actions. Existing Q-guidance methods can be categorized into two classes. The first backpropagates Q-gradients from the final denoised action through all diffusion steps (Wang et al., 2022). While effective in principle, this approach suffers from vanishing or exploding gradients as the number of denoising steps increases, leading to unstable optimization. The second estimates Q-gradients of noisy actions directly at intermediate denoising steps (Fang et al., 2024), thus avoiding backpropagation through the entire trajectory. Although more stable, this method produces unreliable Q-gradients when noisy actions lie far from the data distribution, resulting in suboptimal guidance.

To address the limitations of the second class of methods, we propose QUAD, a Q-gradient uncertainty-aware guidance framework that improves the reliability of denoising guidance. Our

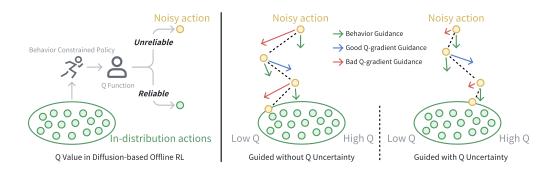


Figure 1: Left: In offline RL, behavior cloning regularization makes the learned Q-function more reliable near the dataset distribution (green), while yielding highly uncertain estimates for out-of-distribution noisy actions (orange). Right: QUAD leverages a Q-ensemble to estimate the uncertainty of Q-gradients and adaptively down-weights unreliable guidance during denoising.

key observation is that critics trained on offline data often yield highly unreliable Q estimates for noisy actions, particularly those far from the dataset distribution (Figure 1, left). To overcome this issue, QUAD employs a Q-ensemble to estimate gradient uncertainty and adaptively attenuate unreliable guidance signals (Figure 1, right). We further provide a theoretical analysis of Q-gradient uncertainty and derive an optimal weighting scheme that minimizes the alignment risk along oracle Q-gradient. Building on this analysis, we design a practical uncertainty-aware weighting mechanism that approximates the theoretical optimum. By integrating this mechanism into the Q-guidance process, QUAD effectively suppresses unreliable gradients, thereby enhancing policy performance.

We evaluate QUAD on the widely adopted D4RL benchmark (Fu et al., 2020), comparing it against state-of-the-art offline RL methods, including both non-diffusion and diffusion-based approaches. Experimental results show that QUAD consistently outperforms prior methods on most tasks and achieves comparable performance on the remaining ones.

In summary, our contributions are threefold:

- We identify and theoretically analyze the limitations of existing Q-guidance methods in diffusion-based offline RL, showing how Q-gradient uncertainty undermines reliability.
- We propose **QUAD**, a novel uncertainty-aware guidance framework that leverages a Q-ensemble to estimate gradient uncertainty and adaptively down-weight unreliable signals.
- We conduct extensive experiments on D4RL, demonstrating that QUAD achieves state-ofthe-art performance across a diverse set of offline RL tasks.

## 2 Preliminaries

A reinforcement learning (RL) problem is typically formulated as a Markov Decision Process (MDP), represented by the tuple  $(\mathcal{S}, \mathcal{A}, \mathcal{T}, r, d_0, \gamma)$ , where  $\mathcal{S}$  denotes the state space,  $\mathcal{A}$  the action space,  $\mathcal{T}(s'|s,a)$  the transition dynamics, r(s,a) the reward function,  $d_0(s)$  the initial state distribution, and  $\gamma \in (0,1)$  the discount factor. The objective of RL is to learn a policy  $\pi(a|s)$  that maximizes the expected discounted cumulative reward (Sutton et al., 1998):

$$J(\pi) = \mathbb{E}_{\pi, \mathcal{T}, d_0} \left[ \sum_{t=0}^{\infty} \gamma^t r(s_t, a_t) \right]$$
 (1)

Offline Reinforcement Learning. Offline RL focuses on learning an effective policy solely from a fixed dataset  $\mathcal{D} = \{(s_i, a_i, r_i, s_i')\}_{i=1}^N$ , which is generated by an (often unknown) behavior policy  $\pi_\beta$ , without access to further environment interactions (Levine et al., 2020). A central challenge in offline RL arises from the distributional shift between  $\pi_\beta$  and the learned policy  $\pi$ , which can lead

to erroneous value estimates. To mitigate this issue, many approaches optimize the expected return under  $Q^{\pi}(s,a)$  while constraining the learned policy to remain close to the behavior policy (Wu et al., 2019):

$$\max_{\pi} \mathbb{E}_{s \sim \mathcal{D}, a \sim \pi(\cdot \mid s)}[Q^{\pi}(s, a)] \quad \text{s.t.} \quad D(\pi \parallel \pi_{\beta}) < \epsilon$$
 (2)

where  $D(\cdot, \cdot)$  denotes a divergence measure (e.g., KL divergence) and  $\epsilon$  is a tolerance parameter.

**Diffusion models.** Diffusion models (Ho et al., 2020; Song et al., 2020a) are a class of generative models that assume latent variables follow a Markovian noising-denoising process. In the forward process  $\{x_{0:T}\}$ , Gaussian noise is gradually added to the clean data  $x_0 \sim p(x_0)$  according to a predefined variance schedule  $\{\beta_{1:T}\}$ :

$$q(x_{1:T}|x_0) = \prod_{t=1}^{T} q(x_t|x_{t-1}), \quad q(x_t|x_{t-1}) := \mathcal{N}(x_t; \sqrt{1-\beta_t}x_{t-1}, \beta_t \mathbf{I})$$
(3)

The marginal distribution admits a closed form:

$$q_t(x_t|x_0) = \mathcal{N}(x_t; \sqrt{\bar{\alpha}_t}x_0, (1 - \bar{\alpha}_t)\mathbf{I}), \quad t \in \{1, \dots, T\}$$
(4)

where  $\alpha_t := 1 - \beta_t$  and  $\bar{\alpha}_t := \prod_{s=1}^t \alpha_s$ . Equivalently, a noisy sample can be reparameterized as

$$x_t = \sqrt{\bar{\alpha}_t} x_0 + \sqrt{1 - \bar{\alpha}_t} \epsilon, \quad \epsilon \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$$
 (5)

Denoising diffusion probabilistic models (DDPMs) (Ho et al., 2020) parameterize the reverse process with Gaussian conditionals  $p_{\theta}(x_{t-1}|x_t) = \mathcal{N}(x_{t-1}; \boldsymbol{\mu}_{\theta}(x_t,t), \boldsymbol{\Sigma}_{\theta}(x_t,t))$ , leading to a generative process:  $p_{\theta}(x_{0:T}) = \mathcal{N}(x_T; \mathbf{0}, \mathbf{I}) \prod_{t=1}^T p_{\theta}(x_{t-1}|x_t)$ . In practice, DDPMs predict the noise  $\boldsymbol{\epsilon}$  in Equation (5) using a neural network  $\boldsymbol{\epsilon}_{\theta}(x_t,t)$  to minimize the evidence lower bound loss:

$$\mathcal{L}(\theta) = \mathbb{E}_{x_0 \sim p(x_0), t \sim \mathcal{U}(1,T), \epsilon \sim \mathcal{N}(\mathbf{0},\mathbf{I})} \left[ \| \epsilon - \epsilon_{\theta}(x_t, t) \|^2 \right]$$
 (6)

**Diffusion-based Offline RL.** Following the DDPM framework, diffusion policies model action generation as a state-conditioned denoising process. Specifically, the noise predictor in (Equation (6)) is replaced with a state-conditional network  $\epsilon_{\theta}(a^t, s, t)$  that predicts actions  $a^0 \in \mathcal{A}$  given the state s, where  $a^t$  denotes the noisy action at denoising step t. This formulation recovers standard behavior cloning (BC) when trained on the dataset  $\mathcal{D}$ . In diffusion-based offline RL, however, pure behavior cloning may fail to exploit Q value information. To address this, Q-function guidance can be incorporated to bias the denoising process toward high-value actions. A straightforward approach, as in Diffusion Q-learning (DQL) (Wang et al., 2022), backpropagates the Q-gradient from the final denoised action  $a^0$  through all denoising steps, leading to the following objective:

$$\arg \min_{\pi_{\theta}} \mathcal{L}(\theta) = \mathbb{E}_{(\boldsymbol{s}, \boldsymbol{a}) \sim \mathcal{D}, t \sim \mathcal{U}(1, T), \boldsymbol{\epsilon} \sim \mathcal{N}(\boldsymbol{0}, \boldsymbol{I})} \Big[ \| \boldsymbol{\epsilon} - \boldsymbol{\epsilon}_{\theta}(\boldsymbol{a}^{t}, \boldsymbol{s}, t) \|^{2} \Big]$$

$$- \eta \cdot \mathbb{E}_{\boldsymbol{s} \sim \mathcal{D}, \boldsymbol{a}^{0} \sim \pi_{\theta}} \big[ Q_{\phi}(\boldsymbol{s}, \boldsymbol{a}^{0}) \big]$$
(7)

where the first term corresponds to the denoising objective, and the second term encourages the policy to generate actions with high Q-values. The coefficient  $\eta$  is a hyperparameter that balances behavior cloning against Q-guidance. An alternative strategy, as in DAC (Fang et al., 2024), directly estimates the Q-gradient of noisy actions at each denoising step, leading to the following objective:

-gradient of noisy actions at each denoising step, leading to the following objective: 
$$\arg\min_{\pi_{\theta}} \mathcal{L}(\theta) = \mathbb{E}_{(\boldsymbol{s},\boldsymbol{a}) \sim \mathcal{D}, t \sim \mathcal{U}(1,T), \boldsymbol{\epsilon} \sim \mathcal{N}(\boldsymbol{0},\boldsymbol{I})} \left[ \eta \cdot \| \boldsymbol{\epsilon} - \boldsymbol{\epsilon}_{\theta}(\boldsymbol{a}^{t},\boldsymbol{s},t) \|^{2} + w(t) \cdot \boldsymbol{\epsilon}_{\theta}(\boldsymbol{a}^{t},\boldsymbol{s},t) \cdot \nabla_{\boldsymbol{a}^{t}} Q_{\phi}(\boldsymbol{s},\boldsymbol{a}^{t}) \right]$$
(8)

where w(t) is a step-dependent weight that controls the influence of Q-gradient guidance across denoising steps. Rather than propagating gradients across the full sequence of denoising steps, DAC-style methods reduce the risk of vanishing or exploding gradients, thereby providing more stable optimization.

# 3 METHODS

We now introduce our proposed method, QUAD, which comprises three main components: (1) a theoretical derivation of the optimal uncertainty-aware weighting scheme for Q-gradient guidance; (2) the formulation and implementation of a Q-gradient uncertainty-aware guidance mechanism; and (3) a practical yet principled procedure for policy extraction. An overview of the QUAD framework is shown in Figure 2.

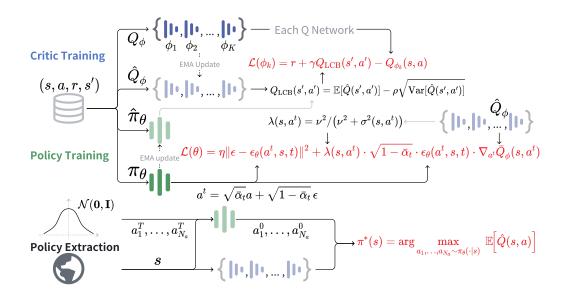


Figure 2: Overview of the QUAD framework: (1) In critic ensemble training, the target policy generates next-step actions and updates critics by minimizing TD error with LCB regularization. (2) In diffusion policy training, the target Q-ensemble estimates Q-gradients and their uncertainty to adaptively down-weight unreliable guidance. (3) In policy extraction, the target diffusion policy proposes candidate actions, and the action with the highest Q-value under the target Q-ensemble is selected.

# 3.1 THEORETICAL ANALYSIS OF Q-GRADIENT UNCERTAINTY-AWARE WEIGHTING

## 3.1.1 UNCERTAINTY IN Q-GRADIENT GUIDANCE

We denote the oracle critic gradient as  $g^*(s, a^t) = \nabla_{a^t} Q^*(s, a^t)$ . The Q-gradient guidance loss in Equation (8) aims to enforce anti-alignment between the noise predictor  $\epsilon_{\theta}(a^t, s, t)$  and  $g^*$ , i.e., to encourage updates along  $-g^*$ . For analytical purposes, let  $u \equiv u_{\theta}(s, a^t, t)$  be the unit vector colinear with  $\epsilon_{\theta}(a^t, s, t)$ . We then introduce the oracle alignment loss:

$$L_{\text{alien}}^* \triangleq \mathbf{u}^\top \mathbf{g}^*. \tag{9}$$

Minimization of  $L_{\text{align}}^*$  promotes anti-alignment with  $g^*$ . Since  $g^*$  is inaccessible in practice, it is replaced by an ensemble-based estimator  $\hat{g}$ :

$$\hat{\boldsymbol{g}}(\boldsymbol{s}, \boldsymbol{a}^t) = \frac{1}{K} \sum_{k=1}^K \nabla_{\boldsymbol{a}^t} Q_{\phi_k}(\boldsymbol{s}, \boldsymbol{a}^t). \tag{10}$$

It is reasonable to assume that  $\hat{g}$  provides an unbiased estimate of  $g^*$ , while exhibiting heteroscedastic variance across different  $(s, a^t)$ . Accordingly,  $\hat{g}$  can be decomposed as the oracle component plus a stochastic perturbation:

$$\hat{\mathbf{g}}(\mathbf{s}, \mathbf{a}^t) = \mathbf{g}^*(\mathbf{s}, \mathbf{a}^t) + \boldsymbol{\xi}(\mathbf{s}, \mathbf{a}^t), \tag{11}$$

with  $\mathbb{E}[\xi] = 0$  and  $\mathrm{Cov}[\xi] = \Sigma(s, a^t)$ . Thus, the alignment loss implemented in practice is

$$L_{\text{align}} \triangleq \boldsymbol{u}^{\top} \hat{\boldsymbol{g}} = \boldsymbol{u}^{\top} \boldsymbol{g}^* + \boldsymbol{u}^{\top} \boldsymbol{\xi}. \tag{12}$$

Here, the perturbation term  $u^{\top} \boldsymbol{\xi}$  has zero expectation and variance  $\sigma^2(\boldsymbol{s}, \boldsymbol{a}^t) = u^{\top} \boldsymbol{\Sigma} u$ , thereby introducing stochastic uncertainty into the alignment objective.

#### 3.1.2 Optimal uncertainty-aware weighting

To mitigate the negative impact of such uncertainty, we introduce a per-sample weight  $\lambda(s, a^t) \in (0, 1]$  to rescale the alignment loss:

$$L_{\text{align}}^{\lambda} \triangleq \lambda \, \boldsymbol{u}^{\top} \hat{\boldsymbol{g}}. \tag{13}$$

Intuitively, samples associated with larger variance should receive smaller weight. Formally, we analyze the problem as one of estimating the oracle alignment  $L^*_{\text{align}}$  from the noisy measurement  $L_{\text{align}}$ . We define the per-sample risk  $\mathcal{R}(\lambda;\sigma^2)$  as the mean-squared error under heteroscedastic noise:

$$\mathcal{R}(\lambda; \sigma^2) \triangleq \mathbb{E}\left[ (L_{\text{align}}^{\lambda} - L_{\text{align}}^*)^2 \mid \sigma^2 \right]$$

$$= \mathbb{E}\left[ (\lambda \mathbf{u}^{\mathsf{T}} \hat{\mathbf{g}} - \mathbf{u}^{\mathsf{T}} \mathbf{g}^*)^2 \mid \sigma^2 \right].$$
(14)

Let  $\nu^2 \triangleq \mathbb{E}[\left(\boldsymbol{u}^{\top}\boldsymbol{g^*}\right)^2]$  denote the second moment of the oracle target. The risk expression can be reformulated as:

$$\mathcal{R}(\lambda; \sigma^2) = (\lambda - 1)^2 \nu^2 + \lambda^2 \sigma^2. \tag{15}$$

The minimizer admits a closed-form solution:

$$\lambda^*(\sigma^2) = \frac{\nu^2}{\nu^2 + \sigma^2}.\tag{16}$$

This solution satisfies  $\lambda^*(0)=1$  and is strictly decreasing in  $\sigma^2$ . Detailed derivations are provided in the Appendix A. Furthermore, Equation (15) explicitly characterizes the bias-variance tradeoff: decreasing  $\lambda$  inflates bias  $(\lambda-1)^2\nu^2$  while reducing variance  $\lambda^2\sigma^2$ . The resulting balance corresponds precisely to *inverse-variance shrinkage*: high-variance samples are systematically down-weighted, whereas noise-free samples retain full weight.

### 3.1.3 A PRACTICAL UNCERTAINTY-AWARE WEIGHTING SCHEME

Although Equation (16) provides the oracle-optimal weight, it cannot be directly computed because  $\nu^2$  is unobserved. Nevertheless, since  $\hat{g}$  is an unbiased estimator of  $g^*$ , the scalar  $u^{\top}\hat{g}$  can be regarded as an unbiased yet heteroscedastic observation of  $u^{\top}g^*$ . This property allows us to estimate  $\nu^2$  from the dataset  $\{(s_i, a_i^t)\}_{i=1}^N$ :

$$\nu^2 \approx \hat{\nu}^2 = \frac{1}{N} \sum_{i=1}^{N} [(\boldsymbol{u}^{\top} \hat{\boldsymbol{g}})^2].$$
 (17)

By substituting this estimate into Equation (16), we obtain a practical weighting scheme:

$$\lambda(\boldsymbol{s}, \boldsymbol{a}^t) = \frac{\hat{\nu}^2}{\hat{\nu}^2 + \sigma^2(\boldsymbol{s}, \boldsymbol{a}^t) + \delta}, \quad \delta > 0.$$
(18)

Here,  $\delta$  is a small positive constant that ensures numerical stability and prevents division by near-zero denominators.

#### 3.2 DIFFUSION POLICY WITH Q-GRADIENT UNCERTAINTY-AWARE GUIDANCE

Building on the above analysis, we implement QUAD by integrating the uncertainty-aware weighting  $\lambda(s, a^t)$  into the diffusion policy training objective Equation (8), following the DAC framework:

$$\mathcal{L}(\theta) = \mathbb{E}_{(\boldsymbol{s}, \boldsymbol{a}^*) \sim \mathcal{D}, t \sim \mathcal{U}(1, T), \epsilon \sim \mathcal{N}(\boldsymbol{0}, \boldsymbol{I})} \Big[ \eta \| \boldsymbol{\epsilon} - \boldsymbol{\epsilon}_{\theta}(\boldsymbol{a}^t, \boldsymbol{s}, t) \|^2 + \lambda(\boldsymbol{s}, \boldsymbol{a}^t) \cdot w(t) \cdot \boldsymbol{\epsilon}_{\theta}(\boldsymbol{a}^t, \boldsymbol{s}, t) \cdot \nabla_{\boldsymbol{a}^t} \bar{Q}_{\phi}(\boldsymbol{s}, \boldsymbol{a}^t) \Big],$$

$$(19)$$

where  $w(t) = \sqrt{1 - \bar{\alpha}_t}$  modulates the strength of Q-guidance according to the noise level, ensuring that the denoised action remains close to the behavior policy in later diffusion steps.

# 3.2.1 Q-Ensemble Training.

A central component of QUAD is the Q-ensemble, which provides unbiased and diverse gradient estimates. To reduce overestimation bias, we train the ensemble using a pessimistic Q-learning scheme (Ghasemipour et al., 2022). Concretely, we maintain K parameterized Q-networks

 $\{Q_{\phi_k}\}_{k=1}^K$  with corresponding target networks  $\{\hat{Q}_{\phi_k}\}_{k=1}^K$ , and employ a lower confidence bound (LCB) as the target value. The critic learning objective for each network is:

$$\mathcal{L}(\phi_i) = \mathbb{E}_{(\boldsymbol{s}, \boldsymbol{a}, \boldsymbol{r}, \boldsymbol{s'}) \sim \mathcal{D}, \, \boldsymbol{a'} \sim \pi_{\theta}} \left[ \boldsymbol{r} + \gamma Q_{\text{LCB}}(\boldsymbol{s'}, \boldsymbol{a'}) - Q_{\phi_i}(\boldsymbol{s}, \boldsymbol{a}) \right]^2,$$

$$Q_{\text{LCB}}(\boldsymbol{s'}, \boldsymbol{a'}) = \mathbb{E}[\hat{Q}(\boldsymbol{s'}, \boldsymbol{a'})] - \rho \sqrt{\text{Var}[\hat{Q}(\boldsymbol{s'}, \boldsymbol{a'})]},$$
(20)

where  $\rho \geq 0$  controls the degree of pessimism, and  $\mathbb{E}[\hat{Q}]$  and  $\text{Var}[\hat{Q}]$  denote the empirical mean and variance across the target critics.

#### 3.2.2 Uncertainty-Aware Weighting.

For stable training, we maintain a target policy  $\hat{\epsilon}_{\theta}$  and target Q-ensemble  $\{\hat{Q}_{\phi_k}\}_{k=1}^K$  using exponential moving averages (EMA). Given a batch  $\{(s,a)\}_1^B$  sampled from  $\mathcal{D}$ , we first add noise to a following Equation (5), yielding  $\{(s,a^t)\}_1^B$ . We then compute the alignment loss of each Q-network relative to the predicted noise direction:

$$l_{align}^{k}(\boldsymbol{s}, \boldsymbol{a}^{t}) = \boldsymbol{u}_{\theta}^{\top} \nabla_{\boldsymbol{a}^{t}} \hat{Q}_{\phi_{k}}(\boldsymbol{s}, \boldsymbol{a}^{t}), \quad \boldsymbol{u}_{\theta} = \frac{\boldsymbol{\epsilon}_{\theta}(\boldsymbol{a}^{t}, \boldsymbol{s}, t)}{\|\boldsymbol{\epsilon}_{\theta}(\boldsymbol{a}^{t}, \boldsymbol{s}, t)\|}.$$
(21)

It is then natural to estimate the heteroscedastic uncertainty  $\sigma^2(s, a^t)$  using the variance of  $\{l_{alian}^k(s, a^t)\}_{k=1}^K$ :

$$\sigma^{2}(\boldsymbol{s}, \boldsymbol{a}^{t}) = \frac{1}{K} \sum_{k=1}^{K} \left( l_{align}^{k}(\boldsymbol{s}, \boldsymbol{a}^{t}) - \bar{l}_{align}(\boldsymbol{s}, \boldsymbol{a}^{t}) \right)^{2}, \quad \bar{l}_{align}(\boldsymbol{s}, \boldsymbol{a}^{t}) = \frac{1}{K} \sum_{k=1}^{K} l_{align}^{k}(\boldsymbol{s}, \boldsymbol{a}^{t}). \quad (22)$$

Following Equation (17), for computational efficiency, we estimate  $\nu^2$  using the empirical variance of  $\{\bar{l}_{align}(s_i, a_i^t)\}_{i=1}^B$  within the batch, and update it using EMA:

$$\hat{\nu}^2 = \frac{1}{B} \sum_{i=1}^{B} (\bar{l}_{align}(\boldsymbol{s}_i, \boldsymbol{a}_i^t) - \bar{\bar{l}}_{align})^2, \quad \bar{\bar{l}}_{align} = \frac{1}{B} \sum_{i=1}^{B} \bar{l}_{align}(\boldsymbol{s}_i, \boldsymbol{a}_i^t). \tag{23}$$

Finally, we compute the uncertainty-aware weight  $\lambda(s, a^t)$  according to Equation (18), and introduce a temperature hyperparameter  $\tau$  to control the aggressiveness of the weighting:

$$\lambda(\mathbf{s}, \mathbf{a}^t) = \frac{\hat{\nu}^2}{\hat{\nu}^2 + \sigma^2(\mathbf{s}, \mathbf{a}^t)/\tau + \delta}.$$
 (24)

When  $\tau \to 0$ , the weighting becomes more aggressive, sharply suppressing high-uncertainty gradients, while as  $\tau \to \infty$ , it approaches uniform weighting. The complete QUAD training procedure is summarized in Algorithm 1.

# 3.3 POLICY EXTRACTION

We denote  $\pi_{\theta}(a|s)$  as the diffusion policy trained via the denoising process with noise predictor  $\epsilon_{\theta}(a^t, s, t)$ . While  $\pi_{\theta}$  can directly generate actions, we further aim to reduce uncertainty during evaluation. To this end, we draw a small batch of  $N_a$  candidate actions from  $\pi_{\theta}(\cdot|s)$  and select the one with the highest ensemble-mean Q-value:

$$\pi^*(s) = \arg \max_{\boldsymbol{a}_1, \dots, \boldsymbol{a}_{N_a} \sim \pi_{\theta}(\cdot | s)} \mathbb{E} \left[ \hat{Q}(s, \boldsymbol{a}) \right]. \tag{25}$$

This extraction strategy is commonly employed in settings where a stochastic actor is used for critic learning, but a deterministic policy is deployed at evaluation. Because  $\pi_{\theta}$  is already trained to approximate the target policy, only a small number of samples  $N_a$  is needed. In our experiments, QUAD achieves strong performance with  $N_a=10$  following DAC, whereas SfBC and Diffusion Q-learning typically require  $N_a=32$  and  $N_a=50$ , respectively.

# 4 RELATED WORK

#### 4.1 OFFLINE RL

Offline RL aims to learn policies from fixed datasets, but suffers from distribution shift that leads to value overestimation in the bootstrapping process (Levine et al., 2020). To address this issue, prior works have developed strategies such as behavior regularization, conservative value estimation, and explicit Bellman error modeling. Behavior-regularized methods constrain policies to stay close to the behavior distribution via candidate-action generators or divergence penalties (Fujimoto et al., 2019; Kumar et al., 2019; Wu et al., 2019). Conservative methods alleviate OOD effects by either adding regularizers to the Q-learning objective (Kumar et al., 2020) or by learning in-sample conservative value functions (Kostrikov et al., 2021; Xu et al., 2023). Alternative approaches explicitly model Bellman errors with a Gumbel distribution and directly learn soft value functions without requiring action sampling (Garg et al., 2023). Our work is related to both behavior-regularized and conservative approaches, as we model the behavior distribution with a diffusion policy and mitigate overestimation bias via a pessimistic Q-ensemble.

#### 4.2 DIFFUSION MODELS

Diffusion models are a class of generative models that consist of a forward diffusion process and a reverse denoising process (Ho et al., 2020), which can also be interpreted as stochastic differential equations (Song et al., 2020b). In the forward process, Gaussian noise is gradually added to the data according to a variance schedule. In the reverse process, a neural network is trained to predict the noise and iteratively recover the clean data. Several works improve efficiency by reducing the number of denoising steps (Song et al., 2020a; Nichol & Dhariwal, 2021; Song et al., 2023). Others explore alternative guidance strategies, such as classifier guidance (Dhariwal & Nichol, 2021) and classifier-free guidance (Ho & Salimans, 2022). More recently, diffusion models have been extended to sequential decision-making, where they are used to represent policies or trajectories (Janner et al., 2022; Chi et al., 2023; Black et al., 2023). Our work builds on diffusion policies and introduces a novel uncertainty-aware Q-gradient guidance mechanism to enhance policy learning in offline RL.

## 4.3 DIFFUSION-BASED OFFLINE RL

Diffusion-based offline RL combines diffusion models with offline RL techniques. A straightforward approach performs behavior cloning with diffusion models and then applies value-based selection to choose high-value actions from the diffusion prior (Chen et al., 2022; Hansen-Estruch et al., 2023). To reduce multi-step sampling cost, an efficient variant distills the diffusion prior into a one-step Gaussian policy (Chen et al., 2024). Another line of work integrates Q-value information directly into diffusion policy training (Wang et al., 2022). However, this approach requires backpropagating Q-gradients through the entire denoising chain, which often causes vanishing or exploding gradients. A more refined strategy applies Q-gradient guidance at each intermediate denoising step, rather than through all steps, as in DAC (Fang et al., 2024). Subsequent extensions incorporate advantage modules or pathwise regularization to further stabilize training (Chen et al., 2025; Gao et al., 2025). Despite their improved stability, these methods still suffer from unreliable Q-gradients when noisy actions deviate from the dataset distribution. Our method addresses this limitation by employing a Q-ensemble to estimate gradient uncertainty and suppress unreliable guidance, thereby improving the robustness of diffusion-based offline RL.

## 5 EXPERIMENTS

In our experiments, we aim to address the following questions:

- Does QUAD outperform state-of-the-art offline RL methods across diverse tasks?
- What is the effect of uncertainty-aware weighting on policy learning and performance?
- How sensitive is QUAD to the choice of uncertainty temperature  $\tau$ ?

Table 1: **Average normalized scores of QUAD vs. baselines.** Abbreviations: "m" = medium, "r" = replay, "e" = expert, "u" = umaze, "div" = diverse, "l" = large. Bold numbers denote the best scores, or the second-best if achieved by our method.

Dataset	Onestep-RL	CQL	IQL	IVR	EQL	Diffuser	DTQL	AlignIQL	SfBC	DQL	DAC	QUAD (ours)
halfcheetah-m	48.4	44.0	47.4	48.3	48.3	44.2	57.9	46.0	45.9	51.1	59.1	<b>62.5</b> ± 0.6
hopper-m	59.6	58.5	66.3	75.5	74.2	58.5	99.6	56.1	57.1	90.5	101.2	$98.7 \pm 2.9$
walker2d-m	81.8	72.5	78.3	84.2	84.2	79.7	89.4	78.5	77.9	87.0	96.8	<b>92.9</b> $\pm$ 6.3
halfcheetah-m-r	38.1	45.5	44.2	44.8	45.2	42.2	50.9	41.1	37.1	47.8	55.0	$57.2 \pm 0.7$
hopper-m-r	97.5	95.0	94.7	99.7	100.7	96.8	100.0	74.8	86.2	101.3	103.1	$104.9 \pm 0.2$
walker2d-m-r	49.5	77.2	73.9	81.2	82.2	61.2	88.5	76.5	65.1	95.5	96.8	<b>99.9</b> ± 0.4
halfcheetah-m-e	93.4	91.6	86.7	94.0	94.2	79.8	92.7	89.1	92.6	96.8	99.1	$100.1 \pm 0.4$
hopper-m-e	103.3	105.4	91.5	111.8	111.2	107.2	109.3	107.1	108.6	111.1	111.7	$111.1 \pm 1.7$
walker2d-m-e	113.0	108.8	109.6	110.2	112.7	108.4	110.0	111.9	109.8	110.1	113.6	$115.5 \pm 1.2$
locomotion total	684.6	698.5	749.7	749.7	752.9	678.0	798.3	681.1	680.3	791.2	836.4	842.8
antmaze-u	64.3	74.0	87.5	93.2	93.8	-	94.8	94.8	92.0	93.4	99.5	$100.0 \pm 0.0$
antmaze-u-div	60.7	84.0	62.2	74.0	82.0	-	78.8	82.4	85.3	66.2	85.0	<b>85.0</b> $\pm$ 4.6
antmaze-m-play	0.3	61.2	71.2	80.2	76.0	-	79.6	80.5	81.3	76.6	85.8	<b>89.5</b> ± 1.7
antmaze-m-div	0.0	53.7	70.0	79.1	73.6	-	82.2	85.5	82.0	78.6	84.0	<b>90.5</b> $\pm$ 4.5
antmaze-l-play	0.0	15.8	39.6	53.2	46.5	-	52.0	65.2	59.3	46.4	50.3	<b>61.0</b> $\pm$ 3.0
antmaze-l-div	0.0	14.9	47.5	52.3	49.0	-	66.4	54.0	45.5	56.6	55.3	$63.0 \pm 5.4$
antmaze total	125.3	303.6	378.0	432.0	420.9	-	441.4	474.8	445.4	417.8	459.9	489.0

## 5.1 SETUP

Offline Datasets. We evaluate QUAD on the widely used D4RL benchmark (Fu et al., 2020), which covers a variety of continuous control tasks with different dataset compositions. Specifically, we consider standard locomotion tasks (HalfCheetah, Hopper, Walker2d) and the more challenging AntMaze tasks. For locomotion, we use version "v0" datasets of three quality levels: medium (m), medium-replay (m-r), and medium-expert (m-e). For AntMaze, we use version "v2" datasets: umaze (u), umaze-diverse (u-div), medium-play (m-play), medium-diverse (m-div), large-play (l-play), and large-diverse (l-div).

**Baselines.** We compare QUAD against a range of offline RL methods, including both non-diffusion and diffusion-based approaches. Non-diffusion baselines include One-step RL (Brandfonbrener et al., 2021), CQL (Kumar et al., 2020), IQL (Kostrikov et al., 2021), IVR (Xu et al., 2023) and EQL (Garg et al., 2023). Diffusion-based baselines include, Diffuser (Janner et al., 2022), SfBC (Chen et al., 2022), Diffusion Q-learning (DQL) (Wang et al., 2022), DTQL (Chen et al., 2024), AlignIQL (He et al., 2024a), and DAC (Fang et al., 2024).

Implementation Details. We implement QUAD on top of the publicly available DAC codebase (Fang et al., 2024). For fair comparison, we adopt the same network architectures and hyperparameters as DAC for both the diffusion policy and the Q-ensemble, unless otherwise specified. We set the ensemble size to K=10 and the temperature to  $\tau=1.0$  for uncertainty weighting, based on preliminary tuning. All models are trained for 2 million gradient steps and evaluated every 20,000 steps using 10 episodes per evaluation. We report the average normalized scores over 4 random seeds for each task, and the final results are averaged over the last 5 evaluations, which typically exhibit stable performance, following the DAC protocol. For baselines, we use the results reported in their respective papers. A complete summary of experimental configurations is provided in Appendix B.

#### 5.2 Main Results

As shown in Table 1, QUAD outperforms most baselines across a variety of tasks, demonstrating the effectiveness of uncertainty-aware weighting and Q-ensemble learning. In particular, QUAD consistently achieves strong performance on the "medium-replay" locomotion datasets and the "medium, large" AntMaze datasets, where broader dataset coverage is available. We hypothesize that in these datasets, the discrepancy in Q-value accuracy between low-uncertainty and high-uncertainty regions is more pronounced, making uncertainty-aware weighting especially beneficial.

#### 5.3 ABLATION STUDIES

Since our method builds on DAC, which can be viewed as the unweighted variant of QUAD, we reproduce DAC results using the same codebase and training protocol as QUAD, indicated by DAC-Rep. We evaluate both methods on locomotion tasks with "medium-replay" datasets and AntMaze "medium" and "large" tasks, where Q-gradient uncertainty is more prevalent. All hyperparameters are kept identical except for those related to uncertainty modeling. As shown in Table 2, QUAD consistently outperforms DAC-Rep across these tasks and exhibits lower variance across random seeds, indicating more stable learning. These gains can be attributed to QUAD's ability to mitigate the negative impact of unreliable Q-gradients. The improvements are particularly pronounced in the challenging AntMaze tasks, where Q-gradients are especially uncertain during denoising.

Table 2: Uncertainty weight ablation on locomotion "medium-replay" datasets and AntMaze "medium"/"large" tasks, comparing QUAD with its unweighted variant DAC-Rep. QUAD achieves higher returns with lower variance, especially on AntMaze where Q-gradients are highly uncertain.

uncertainty weight	walker2d	hopper	halfcheetah		ant	maze	
uncertainty weight	m-r	m-r	m-r	m-p	m-d	l-p	l-d
w/o. (DAC-Rep)	98.1 ± 1.5	$103.4 \pm 0.2$	$55.3 \pm 0.2$	$88.5 \pm 3.0$	$82.5 \pm 17.7$	$41.5 \pm 24.4$	42.5 ± 11.1
w. (QUAD)	99.9 ± 0.4	$\textbf{104.9} \pm \textbf{0.2}$	$\textbf{57.2} \pm \textbf{0.7}$	89.5 ± 1.7	$90.5 \pm 4.5$	$61.0 \pm 3.0$	$63.0 \pm 5.4$

#### 5.4 Sensitivity Analysis

To examine the sensitivity of QUAD to key hyperparameters, we vary the uncertainty temperature  $\tau \in \{0.1, 0.5, 1.0, 10.0, 100\}$ . We present results on "walker2d-medium" and "walker2d-medium-replay" in Figure 3, while full tasks are reported in Appendix B.4. Our findings indicate that QUAD is more sensitive to  $\tau$  in "medium" than in "medium-replay", likely because the former exhibits a narrower data distribution, making uncertainty estimates less reliable than in "medium-replay".

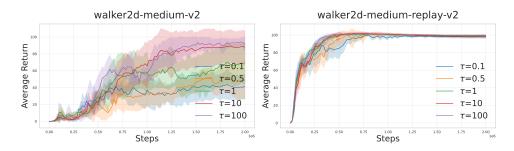


Figure 3: Sensitivity of QUAD to uncertainty temperature  $\tau$ , with stronger effects in "medium" due to its narrower data distribution.

#### 6 Conclusion

We introduced QUAD, a diffusion-based offline RL method that incorporates uncertainty-aware Q-gradient weighting to improve policy learning. By leveraging a Q-ensemble to estimate uncertainty, QUAD mitigates the adverse effects of unreliable Q-gradients during denoising. Our theoretical analysis shows that this weighting scheme stabilizes optimization and enhances policy performance. Extensive experiments on the D4RL benchmark demonstrate that QUAD outperforms state-of-theart diffusion-based methods across diverse tasks, particularly in challenging high-uncertainty settings. A limitation of QUAD lies in its reliance on the variance of Q-ensemble gradients for uncertainty estimation. The diversity of the Q-ensemble is also crucial for reliable uncertainty estimates, which may benefit from techniques such as data augmentation or ensemble diversity promotion. Future work includes exploring more advanced uncertainty estimation methods and extending QUAD to broader RL scenarios, such as offline meta-RL.

# REFERENCES

- Kevin Black, Michael Janner, Yilun Du, Ilya Kostrikov, and Sergey Levine. Training diffusion models with reinforcement learning. *arXiv preprint arXiv:2305.13301*, 2023.
- David Brandfonbrener, Will Whitney, Rajesh Ranganath, and Joan Bruna. Offline rl without off-policy evaluation. *Advances in neural information processing systems*, 34:4933–4946, 2021.
- Huayu Chen, Cheng Lu, Chengyang Ying, Hang Su, and Jun Zhu. Offline reinforcement learning via high-fidelity generative behavior modeling. *arXiv* preprint arXiv:2209.14548, 2022.
- Tianyu Chen, Zhendong Wang, and Mingyuan Zhou. Diffusion policies creating a trust region for offline reinforcement learning. *Advances in Neural Information Processing Systems*, 37:50098–50125, 2024.
- Xuyang Chen, Keyu Yan, and Lin Zhao. Taming ood actions for offline reinforcement learning: An advantage-based approach. *arXiv preprint arXiv:2505.05126*, 2025.
- Cheng Chi, Zhenjia Xu, Siyuan Feng, Eric Cousineau, Yilun Du, Benjamin Burchfiel, Russ Tedrake, and Shuran Song. Diffusion policy: Visuomotor policy learning via action diffusion. *The International Journal of Robotics Research*, pp. 02783649241273668, 2023.
- Prafulla Dhariwal and Alexander Nichol. Diffusion models beat gans on image synthesis. *Advances in neural information processing systems*, 34:8780–8794, 2021.
- Linjiajie Fang, Ruoxue Liu, Jing Zhang, Wenjia Wang, and Bing-Yi Jing. Diffusion actor-critic: Formulating constrained policy iteration as diffusion noise regression for offline reinforcement learning. *arXiv preprint arXiv:2405.20555*, 2024.
- 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.
- Scott Fujimoto and Shixiang Shane Gu. A minimalist approach to offline reinforcement learning. *Advances in neural information processing systems*, 34:20132–20145, 2021.
- Scott Fujimoto, Edoardo Conti, Mohammad Ghavamzadeh, and Joelle Pineau. Benchmarking batch deep reinforcement learning algorithms. *arXiv preprint arXiv:1910.01708*, 2019.
- Chen-Xiao Gao, Chenyang Wu, Mingjun Cao, Chenjun Xiao, Yang Yu, and Zongzhang Zhang. Behavior-regularized diffusion policy optimization for offline reinforcement learning. *arXiv* preprint arXiv:2502.04778, 2025.
- Divyansh Garg, Joey Hejna, Matthieu Geist, and Stefano Ermon. Extreme q-learning: Maxent rl without entropy. *arXiv preprint arXiv:2301.02328*, 2023.
- Kamyar Ghasemipour, Shixiang Shane Gu, and Ofir Nachum. Why so pessimistic? estimating uncertainties for offline rl through ensembles, and why their independence matters. *Advances in Neural Information Processing Systems*, 35:18267–18281, 2022.
- 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.
- Longxiang He, Li Shen, Junbo Tan, and Xueqian Wang. Aligniql: Policy alignment in implicit q-learning through constrained optimization. *arXiv preprint arXiv:2405.18187*, 2024a.
- Tairan He, Zhengyi Luo, Xialin He, Wenli Xiao, Chong Zhang, Weinan Zhang, Kris Kitani, Changliu Liu, and Guanya Shi. Omnih2o: Universal and dexterous human-to-humanoid whole-body teleoperation and learning. *arXiv* preprint arXiv:2406.08858, 2024b.
- Jonathan Ho and Tim Salimans. Classifier-free diffusion guidance. *arXiv preprint* arXiv:2207.12598, 2022.
  - Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. *Advances in neural information processing systems*, 33:6840–6851, 2020.

- Michael Janner, Yilun Du, Joshua B Tenenbaum, and Sergey Levine. Planning with diffusion for flexible behavior synthesis. arXiv preprint arXiv:2205.09991, 2022.
- Ilya Kostrikov. JAXRL: Implementations of Reinforcement Learning algorithms in JAX, 10 2021. URL https://github.com/ikostrikov/jaxrl.
  - Ilya Kostrikov, Ashvin Nair, and Sergey Levine. Offline reinforcement learning with implicit q-learning. *arXiv preprint arXiv:2110.06169*, 2021.
- Aviral Kumar, Justin Fu, Matthew Soh, George Tucker, and Sergey Levine. Stabilizing off-policy q-learning via bootstrapping error reduction. *Advances in neural information processing systems*, 32, 2019.
  - Aviral Kumar, Aurick Zhou, George Tucker, and Sergey Levine. Conservative q-learning for offline reinforcement learning. *Advances in neural information processing systems*, 33:1179–1191, 2020.
  - Guillaume Lample and Devendra Singh Chaplot. Playing fps games with deep reinforcement learning. In *Proceedings of the AAAI conference on artificial intelligence*, volume 31, 2017.
  - Sergey Levine, Aviral Kumar, George Tucker, and Justin Fu. Offline reinforcement learning: Tutorial, review, and perspectives on open problems. *arXiv* preprint arXiv:2005.01643, 2020.
  - Misra D Mish. A self regularized non-monotonic activation function. 2019. *arXiv preprint arXiv:1908.08681*, 1908.
  - Volodymyr Mnih, Koray Kavukcuoglu, David Silver, Alex Graves, Ioannis Antonoglou, Daan Wierstra, and Martin Riedmiller. Playing atari with deep reinforcement learning. *arXiv preprint arXiv:1312.5602*, 2013.
  - Alexander Quinn Nichol and Prafulla Dhariwal. Improved denoising diffusion probabilistic models. In *International conference on machine learning*, pp. 8162–8171. PMLR, 2021.
  - Jiaming Song, Chenlin Meng, and Stefano Ermon. Denoising diffusion implicit models. *arXiv* preprint arXiv:2010.02502, 2020a.
  - Yang Song, Jascha Sohl-Dickstein, Diederik P Kingma, Abhishek Kumar, Stefano Ermon, and Ben Poole. Score-based generative modeling through stochastic differential equations. *arXiv* preprint *arXiv*:2011.13456, 2020b.
  - Yang Song, Prafulla Dhariwal, Mark Chen, and Ilya Sutskever. Consistency models. 2023.
  - Richard S Sutton, Andrew G Barto, et al. *Reinforcement learning: An introduction*, volume 1. MIT press Cambridge, 1998.
  - Zhendong Wang, Jonathan J Hunt, and Mingyuan Zhou. Diffusion policies as an expressive policy class for offline reinforcement learning. *arXiv* preprint arXiv:2208.06193, 2022.
  - Yifan Wu, George Tucker, and Ofir Nachum. Behavior regularized offline reinforcement learning. *arXiv preprint arXiv:1911.11361*, 2019.
  - Haoran Xu, Li Jiang, Jianxiong Li, Zhuoran Yang, Zhaoran Wang, Victor Wai Kin Chan, and Xianyuan Zhan. Offline rl with no ood actions: In-sample learning via implicit value regularization. arXiv preprint arXiv:2303.15810, 2023.
  - Yanjie Ze, Zixuan Chen, JoÃÇo Pedro AraÚjo, Zi-ang Cao, Xue Bin Peng, Jiajun Wu, and C Karen Liu. Twist: Teleoperated whole-body imitation system. *arXiv* preprint arXiv:2505.02833, 2025.
  - Hongtu Zhou, Ruiling Yang, Yakun Zhu, Haoqi Zhao, Hai Zhang, Di Zhang, Junqiao Zhao, Chen Ye, and Changjun Jiang. Certain: Context uncertainty-aware one-shot adaptation for context-based offline meta reinforcement learning. In *International conference on machine learning*, 2025.

THE USE OF LLMS

We thank ChatGPT-5 for its assistance in polishing the writing and proofreading of this paper. The authors are responsible for the content and presentation.

# A DETAILED PROOFS AND EXTENSIONS OF QUAD THEORY

In this appendix, we give complete derivations for the risk decomposition and the optimal per-sample weight. Recall that for a state-action pair  $(s, a^t)$  and a unit direction  $u \equiv u_{\theta}(s, a^t, t)$ , the oracle alignment loss and the noisy measurement are

$$L_{\text{align}}^* \triangleq \boldsymbol{u}^\top \boldsymbol{g}^*, \qquad L_{\text{align}} \triangleq \boldsymbol{u}^\top \hat{\boldsymbol{g}} = \boldsymbol{u}^\top \boldsymbol{g}^* + \boldsymbol{u}^\top \xi.$$
 (26)

The rescaled alignment loss is

$$L_{\text{align}}^{\lambda} \triangleq \lambda \, \boldsymbol{u}^{\top} \hat{\boldsymbol{g}}. \tag{27}$$

Let  $\nu^2 \triangleq \mathbb{E}[(u^{\top}g^*)^2]$  denote the second moment of the oracle target along u.

$$\mathcal{R}(\lambda; \sigma^{2}) \triangleq \mathbb{E}\left[\left(L_{\text{align}}^{\lambda} - L_{\text{align}}^{*}\right)^{2} \mid \sigma^{2}\right] \\
= \mathbb{E}\left[\left(\lambda \mathbf{u}^{\top} \hat{\mathbf{g}} - \mathbf{u}^{\top} \mathbf{g}^{*}\right)^{2} \mid \sigma^{2}\right] \\
= \mathbb{E}\left[\left(\lambda \left(\mathbf{u}^{\top} \hat{\mathbf{g}} - \mathbf{u}^{\top} \mathbf{g}^{*}\right) + (\lambda - 1) \mathbf{u}^{\top} \mathbf{g}^{*}\right)^{2} \mid \sigma^{2}\right] \\
= \mathbb{E}\left[\lambda^{2} \left(\mathbf{u}^{\top} \hat{\mathbf{g}} - \mathbf{u}^{\top} \mathbf{g}^{*}\right)^{2} + 2\lambda(\lambda - 1) \left(\mathbf{u}^{\top} \hat{\mathbf{g}} - \mathbf{u}^{\top} \mathbf{g}^{*}\right) \mathbf{u}^{\top} \mathbf{g}^{*} + (\lambda - 1)^{2} \left(\mathbf{u}^{\top} \mathbf{g}^{*}\right)^{2} \mid \sigma^{2}\right] \\
= \lambda^{2} \mathbb{E}\left[\left(\mathbf{u}^{\top} \hat{\mathbf{g}} - \mathbf{u}^{\top} \mathbf{g}^{*}\right)^{2} \mid \sigma^{2}\right] + 2\lambda(\lambda - 1) \mathbb{E}\left[\left(\mathbf{u}^{\top} \hat{\mathbf{g}} - \mathbf{u}^{\top} \mathbf{g}^{*}\right) \mathbf{u}^{\top} \mathbf{g}^{*} \mid \sigma^{2}\right] \\
+ (\lambda - 1)^{2} \mathbb{E}\left[\left(\mathbf{u}^{\top} \mathbf{g}^{*}\right)^{2} \mid \sigma^{2}\right] \\
= \lambda^{2} \mathbb{E}\left[\left(\mathbf{u}^{\top} \boldsymbol{\xi}\right)^{2} \mid \sigma^{2}\right] + 2\lambda(\lambda - 1) \mathbb{E}\left[\left(\mathbf{u}^{\top} \boldsymbol{\xi}\right) \left(\mathbf{u}^{\top} \mathbf{g}^{*}\right) \mid \sigma^{2}\right] + (\lambda - 1)^{2} \mathbb{E}\left[\left(\mathbf{u}^{\top} \mathbf{g}^{*}\right)^{2} \mid \sigma^{2}\right] \\
= \lambda^{2} \sigma^{2} + (\lambda - 1)^{2} \nu^{2} \tag{28}$$

The last equality holds because the perturbation term  $u^{\top} \xi$  has zero expectation and variance  $\sigma^2(s, a^t)$ , while the oracle component  $g^*(s, a^t)$  and the stochastic perturbation  $\xi(s, a^t)$  are independent.

Expanding the quadratic form yields

$$\mathcal{R}(\lambda; \sigma^2) = \lambda^2 \sigma^2 + (\lambda^2 - 2\lambda + 1)\nu^2 = (\nu^2 + \sigma^2)\lambda^2 - 2\nu^2\lambda + \nu^2.$$
 (29)

Since  $\nu^2 + \sigma^2 > 0$ , the function is strictly convex in  $\lambda$ . Taking the derivative and setting it to zero,

$$\frac{\partial \mathcal{R}}{\partial \lambda} = 2(\nu^2 + \sigma^2)\lambda - 2\nu^2 = 0,\tag{30}$$

we obtain the unique minimizer

$$\lambda^*(\sigma^2) = \frac{\nu^2}{\nu^2 + \sigma^2}.\tag{31}$$

# B DETAILS OF EXPERIMENTAL SETUP

We train all models for 2M gradient steps. Each environment is run with 4 independent seeds, and performance is evaluated every 20k steps using 10 additional seeds, yielding 40 rollouts per evaluation. We report the mean score over the final 50k steps without early stopping. Experiments are conducted on 4 RTX 4090 GPUs, with each run taking about 2.5 hours including training and evaluation. Our implementation builds on the jaxrl (Kostrikov, 2021) codebase.

#### **B.1** Network Architecture

Both the actor and critic adopt a 3-layer MLP with hidden size 256 and Mish activation (Mish, 1908). Target networks are used to stabilize training:  $\hat{\epsilon}_{\theta}$  and  $\hat{Q}_{\phi_k}$  are initialized with the same parameters as  $\epsilon_{\theta}$  and  $Q_{\phi_k}$ , and track their exponential moving averages (EMA). The target actor is updated every 5 gradient steps, while the target critics are updated after each step.

### **B.2** Hyperparameters

We use consistent hyperparameter settings for the diffusion models and networks across all tasks. The hyperparameters are specified as follows:

Table 3: Hyperparameters for all networks and tasks.

Hyperparameter	Value			
T (Diffusion Steps)	5			
$\beta_t$ (Noise Schedule)	Variance Preserving			
K (Ensemble Size)	10			
B (Batch Size)	256			
Learning Rates (for all networks)	3e-4, 1e-3 (antmaze-large)			
Learning Rate Decay	Cosine			
Optimizer	Adam			
$\eta_{\text{init}}$ (Initial Behavior Cloning Strength)	[0.1, 1]			
$\alpha_{\eta}$ (for Dual Gradient Ascent)	0.001			
$\alpha_{\rm ema}$ (EMA Learning Rate)	5e-3			
$N_a$ (Number of sampled actions for evaluation)	10			
b (Behavior Cloning Threshold)	[0.05, 1]			
$\rho$ (Pessimistic factor)	[0, 2]			

QUAD adopts the same hyperparameters as DAC for the diffusion policy and Q-ensemble, except for AntMaze "large" tasks where a smaller  $\eta$  is used. We sweep over  $\tau \in \{0.1, 0.5, 1.0, 10.0, 100\}$  and report the best value for each task in Table 4.

Table 4: Hyperparameters settings for tasks.

Tasks	au	b	$\eta$	ρ	Regularization Type
hopper-medium-v2	10	1	-	1.5	Learnable
hopper-medium-replay-v2	0.1	1	-	1.5	Learnable
hopper-medium-expert-v2	0.1	0.05	-	1.5	Learnable
walker2d-medium-v2	100	1	-	1	Learnable
walker2d-medium-replay-v2	100	1	-	1	Learnable
walker2d-medium-expert-v2	100	1	-	1	Learnable
halfcheetah-medium-v2	0.1	1	-	0	Learnable
halfcheetah-medium-replay-v2	1.0	1	-	0	Learnable
halfcheetah-medium-expert-v2	100	0.1	-	0	Learnable
antmaze-umaze-v0	10	-	0.1	1	Constant
antmaze-umaze-diverse-v0	10	-	0.1	1	Constant
antmaze-medium-play-v0	10	-	0.1	1	Constant
antmaze-medium-diverse-v0	10	-	0.1	1	Constant
antmaze-large-play-v0	1	-	0.01	1.1	Constant
antmaze-large-diverse-v0	1	-	0.01	1	Constant

# B.3 PSEUDO CODE OF QUAD

We provide the pseudo code of QUAD in Algorithm 1.

```
702
             Algorithm 1 QUAD: Q-gradient Uncertainty-aware Guidance Training
703
             Require: offline dataset \mathcal{D}, batch size B, learning rates \alpha_{\phi}, \alpha_{\theta}, \alpha_{\eta} and \alpha_{\text{ema}}, behavior cloning
704
                    threshold \varepsilon_b, pessimism factor \rho, initial Lagrangian multiplier \eta_{\text{init}}, ensemble size K, uncertainty
705
                    temperature \tau
706
               1: Initialize: diffusion policy \epsilon_{\theta}, target diffusion policy \hat{\epsilon}_{\theta} = \epsilon_{\theta}, Q ensemble networks Q_{\phi_k}, target
707
                    Q ensemble networks \hat{Q}_{\phi_k}=Q_{\phi_k} (i=1,2,...,K), Lagrangian multiplier \eta=\eta_{\rm init}
708
               2: while training not convergent do
709
               3:
                         Sample a batch of B transitions \{(s, a, r, s')\} \subset \mathcal{D}
710
                         Sample a' = a^0 through denoising process using noise predictor \hat{\epsilon}_{\theta}(a^t, s, t).
               4:
               5:
                         for k in \{1, 2, ..., K\} do
711
                               Update \phi_k \leftarrow \phi_k - \alpha_\phi \nabla_{\phi_k} \mathcal{L}(\phi_k) (Equation (20))
                                                                                                                                    ▷ Q ensemble learning
               6:
712
               7:
713
                         Sample \epsilon \sim \mathcal{N}(\mathbf{0}, \mathbf{I}), t \sim \mathcal{U}(0, T) and compute \mathbf{a}^t = \sqrt{\bar{\alpha}_t} \mathbf{a} + \sqrt{1 - \bar{\alpha}_t} \epsilon
               8:
714
               9:
                         Estimate Q-gradient \nabla_{\boldsymbol{a^t}} Q_{\pi_i}(\boldsymbol{s}, \boldsymbol{a^t}) using (Equation (10))
715
                         Estimate Q-gradient uncertainty weight \lambda(s, a^t) using (Equation (18))
             10:
716
                         \theta \leftarrow \theta - \alpha_{\theta} \nabla_{\theta} \mathcal{L}(\theta) (Equation (19))
                                                                                                                                             ⊳ Policy learning
             11:
717
                         \eta \leftarrow \eta + \alpha_{\eta}(||\boldsymbol{\epsilon}_{\theta}(\boldsymbol{a^t}, \boldsymbol{s}, t) - \boldsymbol{\epsilon}||^2 - \varepsilon_b)
                                                                                                                    Dual gradient ascent (optional)
             12:
718
                         \hat{\theta} \leftarrow (1 - \alpha_{\text{ema}})\hat{\theta} + \alpha_{\text{ema}}\theta
             13:
719
                         \hat{\phi}_i \leftarrow (1 - \alpha_{\text{ema}})\phi_i + \alpha_{\text{ema}}\phi_i
             14:
                                                                                                             ▶ Update target networks using EMA
720
             15: end while
721
```

## **B.4** SENSITIVITY ANALYSIS

722 723

724 725

726

We present the full sensitivity analysis of QUAD to uncertainty temperature  $\tau$  on all tasks in Figure 4.

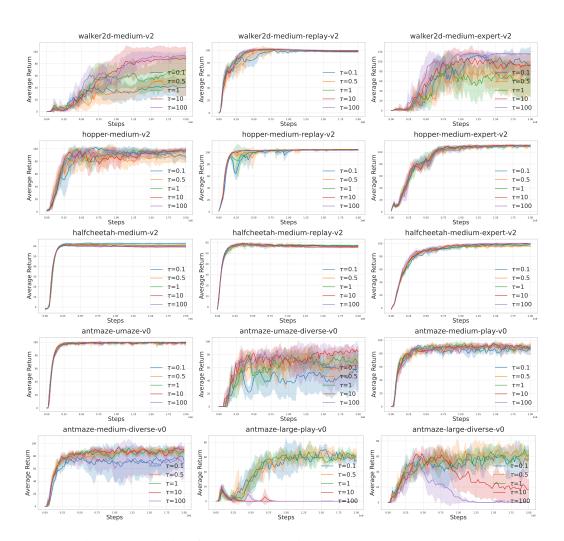


Figure 4: Sensitivity of QUAD to uncertainty temperature  $\tau$  on various tasks.