Q-SUPERVISED CONTRASTIVE REPRESENTATION: A STATE DECOUPLING FRAMEWORK FOR SAFE OFFLINE REINFORCEMENT LEARNING

Anonymous authors

Paper under double-blind review

ABSTRACT

Safe offline reinforcement learning (RL), which aims to learn the safety-guaranteed policy without risky online interaction with environments, has attracted growing recent attention for safety-critical scenarios. However, existing approaches encounter out-of-distribution problems during the testing phase, which can result in potentially unsafe outcomes. This issue arises due to the infinite possible combinations of reward-related and cost-related states. In this work, we propose State Decoupling with Q-supervised Contrastive representation (SDQC), a novel framework that decouples the global observations into reward- and cost-related representations for decision-making, thereby improving the generalization capability for unfamiliar global observations. Compared with the classical representation learning methods, which typically require model-based estimation (e.g., bisimulation), we theoretically prove that our Q-supervised method generates a coarser representation while preserving the optimal policy, resulting in improved generalization performance. Experiments on DSRL benchmark problems provide compelling evidence that SDOC surpasses other baseline algorithms, especially for its exceptional ability to achieve almost zero violations in more than half of the tasks, while the state-of-theart algorithm can only achieve the same level of success in a quarter of the tasks. Further, we demonstrate that SDQC possesses superior generalization ability when confronted with unseen environments.

031 032

033 034

000

001

002

004

006

008 009 010

011 012 013

014

015

016

017

018

019

021

024

025

026

027

028

029

1 INTRODUCTION

Reinforcement learning (RL) has been proven to be a powerful tool for solving high-dimensional decision-making problems under uncertainty (Mnih et al., 2015; Silver et al., 2017; Schrittwieser et al., 2020). Nevertheless, safety concerns remain a significant obstacle to the extensive adoption of RL in safety-critical domains (Garcia & Fernández, 2015; Gu et al., 2022; Xu et al., 2022b; Li, 2023), such as industrial management, and robot control. In these contexts, the potential for catastrophic outcomes necessitates a significant emphasis on preventing unsafe actions (Andersen et al., 2020; Kiran et al., 2021; Brunke et al., 2022). As a promising method that received growing attention, safe RL provides safety guarantees by formulating the problem as a constrained Markov decision process (CMDP) (Altman, 1998; 2021).

043 Over the past few years, a multitude of safe RL algorithms have been introduced (Achiam et al., 044 2017; Tessler et al., 2018; Chow et al., 2018; Zhao et al., 2021; Sootla et al., 2022; Yu et al., 2022). Regrettably, most existing methodologies address safety concerns within online settings. Such online 046 methods rely on the availability of high-fidelity simulators or involve direct agent-environment 047 interactions during the training process, thereby introducing additional risks of safety violations (Liu 048 et al., 2023a). Safe offline RL, on the other hand, provides a promising solution that learns the safetyguaranteed policy in a fully offline manner. Its training requires no additional risky interaction with the environment and relies only on the pre-collected offline dataset. However, empirical observations 051 indicate that most existing safe offline RL algorithms fail to thoroughly ensure pre-defined safety constraints during testing (Liu et al., 2023a; Zheng et al., 2024). Such occurrences tend to be more 052 pronounced in tasks characterized by higher complexity or in environments with higher observation dimensions.

In offline RL, it is imperative that the states vis-055 ited during testing have been included in (or at least 056 not far away from) the training dataset to ensure 057 robust performance (Fujimoto et al., 2019; Wang 058 et al., 2022). However, Safe offline RL problems have various combinations of reward-related and cost-related states. For instance, as illustrated in 060 Figure 1, a UGV (unmanned ground vehicle) needs 061 to navigate around traps to reach its final destination. 062 During testing, if the relative positions of traps and 063 the target haven't occurred in the training dataset, 064 the agent may struggle to make informed decisions 065 based on such unknown observations. It is reason-066 able to suspect that the primary reason for the subpar 067 performance of safe offline RL during tests lies in 068 the out-of-distribution (OOD) issue.



Figure 1: Diagram of the OOD issue for offline trained UGV in the testing phase.

To tackle this problem and improve the generaliza-

tion of safe offline RL, we propose State Decoupling with Q-supervised Contrastive representation 071 (SDQC), a novel representation learning method that decouples the global observations into reward-072 and cost-related representations. Attributable to the successful application of Hamilton-Jacobi (HJ) 073 reachability analysis in Safe RL (Fisac et al., 2019; Yu et al., 2022; Zheng et al., 2024), which 074 introduces a safety analysis method iterated through Q-learning with convergence guarantees, our 075 approach conducts safety assessments on the cost-related representations and make decisions based on the assessment results. Our SDQC, developed based on FISOR (Zheng et al., 2024), distinguishes 076 itself from FISOR and other classical methods which rely on global observations for decision-making 077 (as depicted in the left subplot of Figure 2), by being the **first** to utilize decoupled representations for decision-making in safe RL tasks (see the right subplot of Figure 2). It employs reward-related 079 representations to make decisions when the assessment confirms absolute safety, switches to costrelated representations when the assessment deems the situation unsafe, and integrates both when the 081 assessment indicates borderline safety. 082

Nevertheless, effective differentiation between reward- and cost-related information from global observations poses a formidable challenge, especially when certain dimensions of the observations contain intertwined information. For instance, some information, like speed and acceleration of UGV, should be included in both reward- and cost-related representations. On the other hand, information from the environment-detecting sensors, which include positions of destinations and obstacles, should be distinctly decoupled. Manual separation of such information proves impractical in most cases.

Towards this end, our Q-supervised contrastive representation decouples the global observations 089 through clustering representations that demonstrate similar learned-Q* across the actions in support[§]. 090 The representations solely capture either reward or cost information, independent of another factor, 091 as determined by the training of Q*. In contrast to model-based representations learning (e.g., 092 bisimulation), our SDQC circumvents the need for model estimation, thus mitigating the challenges 093 posed by severe estimation errors in scenarios with sparse rewards or costs. Moreover, we demonstrate 094 that our representations can be trained concurrently with the Q*-learning process by incorporating an 095 additional loss term within the framework of implicit Q-learning. 096

Further, we provide theoretical evidence that our method produces a coarser representation compared to bisimulation, while still preserving the optimal policy. This is supported by our argument that SDQC leads to a higher information entropy of the global observations when conditioned on the representations. This attribute grants SDQC superior generalization capabilities, bolstering its efficacy in handling OOD observations during the testing phase.

The experimental results showcase that our SDQC outperforms other safe offline RL algorithms in the DSRL benchmark, especially in its exceptional ability to achieve zero violations in the majority of tasks. Further, in generalization tests where agents are evaluated in environments with a different number of obstacles than those in the training dataset, all baseline algorithms show a substantial

106 107

[§]For simplicity, we denote Q^* as a generic notation to represent the optimal Q-value functions for both reward (Q_r^*) and $\cot(Q_h^*)$ throughout this paper. Similarly, Q refers to both Q_r and Q_h .



Figure 2: Overview diagram of classical approaches for safe decision-making (left) and our proposed state-decoupling framework for safe decision-making (right).

increase in cost and/or a significant decline in reward. In contrast, SDQC stands out as the only approach that guarantees no increase in cost while experiencing only a slight decline in reward.

2 PRELIMINARIES

Safe Offline RL. Safe RL tasks are generally modeled as CMDP in the form of \mathcal{M} = 132 $(\mathcal{S}, \mathcal{A}, P, r, c, \gamma, d_0)$, where \mathcal{S} is the state space, \mathcal{A} is the action space, P is the model dynam-133 ics, $r: S \times A \to \mathbb{R}$ represents the reward function, $c: S \times A \to [0, C_{max}]$ represents the cost 134 function, $\gamma \in [0,1)$ is the discount factor, and $d_0 \in \Delta(S)$ is the distribution of initial state s_0 135 (the set of the probability distribution over S is denoted as $\Delta(S)$). $P(s'|s,a) : S \times A \to \Delta(S)$ 136 represents the transition function from state s to s' when taking action a. The state-action-reward-cost 137 transitions over trajectory are recorded as $\tau := (s_t, a_t, r_t, c_t)_{t>0}$. The goal of Safe RL is to learn a 138 policy $\pi: S \to \Delta(A)$ that maximizes the expectation of the cumulated discounted reward while 139 restricting the expected cumulative costs below a predefined cost limit κ , which can be denoted by $\max_{\pi} \mathbb{E}_{\tau \sim \pi} [\sum_{t=0}^{\infty} \gamma^t r(s_t, a_t)]$, s.t. $\mathbb{E}_{\tau \sim \pi} [\sum_{t=0}^{\infty} \gamma^t c(s_t, a_t)] \leq \kappa$. 140 141

In offline settings, the training is performed on a statistical dataset denoted as $\mathcal{D}_{\beta} := (s, a, s', r, c)$. This offline dataset comprises both safe and unsafe trajectories and is acquired from a behavior policy π_{β} . During training, most existing safe offline RL algorithms utilize the temporal difference (TD) method to learn the reward state-value function $V_r^{\pi}(s_t)$, which models the expected cumulative reward $\mathbb{E}_{\tau \sim \pi}[\sum_{i=t}^{\infty} \gamma^i r(s_i, a_i)]$, as well as the cost state-value function $V_c^{\pi}(s_t)$, which models the expected cumulative cost $\mathbb{E}_{\tau \sim \pi}[\sum_{i=t}^{\infty} \gamma^i c(s_i, a_i)]$. The primal training objective of safe offline RL can be expressed as follows:

148 149

125

126 127

128

129 130

131

$$\max_{\pi} \mathbb{E}_{s_t \sim \mathcal{D}_{\beta}}[V_r^{\pi}(s_t)] \qquad \text{s.t.} \ \mathbb{E}_{s_t \sim \mathcal{D}_{\beta}}[V_c^{\pi}(s_t)] \le \kappa; \ D(\pi | \pi_{\beta}) \le \epsilon_{\pi}, \tag{1}$$

150 151

where $D(\pi | \pi_{\beta})$ is the divergence term that prevents the distributional shift in offline training.

152 A commonly employed approach for solving Eq. 1 involves reformulating the training objective using the Lagrangian dual form as $\min_{\lambda>0} \max_{\pi} \mathbb{E}_{s_t \sim \mathcal{D}_{\beta}}[V_r^{\pi}(s_t) - \lambda(V_c^{\pi}(s_t) - \kappa)]$, s.t. $D(\pi|\pi_{\beta}) \leq \epsilon_{\pi}$, 153 where the learnable Lagrange multiplier λ is iteratively updated to enforce the constraint. However, the 154 Lagrangian approach can be sensitive to the learning rate and initialization of the multiplier (Stooke 155 et al., 2020). Furthermore, the joint optimization of V_r^{π} , V_c^{π} , and π leads to significant instability, 156 as minor approximation errors can bootstrap across them and propagate, thereby undermining the 157 ability to provide robust safety guarantees (Kumar et al., 2019; Zheng et al., 2024). In safety-critical 158 scenarios, where violations are deemed unacceptable, there is a pressing need for advancements in 159 addressing these challenges. 160

Hamilton-Jacobi reachability. As a prospective method to perform safety assessment rooted in control theory, HJ reachability has been proven to be applicable in Safe RL tasks for both online

162 settings (Chen et al., 2021a; Yu et al., 2022) and offline settings (Zheng et al., 2024). In addition to 163 the tuple formulation \mathcal{M} in CMDP, we introduce a constraint violation function $h: \mathcal{S} \to \mathbb{R}$, which is 164 positive if the state constraint is violated and negative otherwise. For a given state s, the safe value 165 function $V_h^{\pi}(s) := \max_{t \in \mathbb{N}} \{h(s_t) \mid s_0 = s, a_i \sim \pi(\cdot | s_i), \forall i \in \{0, \dots, t\}\}$ represents the worst constraint violations among all possible trajectories induced by policy π . The corresponding safe 166 Q-function is given by $Q_h^{\pi}(s, a) := \max_{t \in \mathbb{N}} \{h(s_t) \mid s_0 = s, a_0 = a, a_i \sim \pi(\cdot | s_i), \forall i \in \{1, \ldots, t\}\}.$ 167 The optimal safe value function, defined as $V_h^*(s) := \min_{\pi} V_h^{\pi}(s)$, stands for the smallest violation 168 one can obtain through adjusting the policy π . Similarly, the corresponding optimal safe Q-function can be expressed as $Q_h^*(s, a) := \min_{\pi} Q_h^{\pi}(s, a)$. With the discount factor defined as γ , Fisac et al. 170 (2019) introduce the following safety Bellman operators: 171

$$\mathcal{B}_{h}^{*}Q_{h,\gamma}(s,a) := (1-\gamma)h(s) + \gamma \max\{h(s), V_{h,\gamma}(s')\}, \quad V_{h,\gamma}(s') = \min_{a'} Q_{h}(s',a'), \quad (2)$$

which is a contraction mapping satisfying $\lim_{\gamma \to 1} Q_{h,\gamma} \to Q_h^*$, $\lim_{\gamma \to 1} V_{h,\gamma} \to V_h^*$. A direct safety inference can be made after training converges. $V_h^*(s) \le 0$ implies the existence of policies that guarantee adherence to the hard constraints throughout the trajectory. Conversely, $V_h^*(s) > 0$ indicates that the destiny towards unsafe states regardless of the chosen policy. In the offline settings, Zheng et al. (2024) pioneered the application of HJ reachability analysis for safety assessment. They present that the decision-making for Safe RL with hard constraints can be decoupled as:

where $A_r^*(s, a) := Q_r^*(s, a) - V_r^*(s)$ and $A_h^*(s, a) := Q_h^*(s, a) - V_h^*(s)$. Eq. 3 theoretically ensures zero constraint violations. However, challenges arise from estimation errors and OOD problems during the testing phase. As a result, the empirical results demonstrate the inability of their algorithm (FISOR) to achieve absolute safety guarantees.

3 STATE DECOUPLING WITH Q-SUPERVISED CONTRASTIVE REPRESENTATIONS

190 191 192

193

194

195 196

197

207

172

179

181 182 183

184

185

186

187 188

189

In our state-decoupling framework, we aim to decouple the state s into two separate representations: one related to rewards, denoted as s_r , and the other related to costs, denoted as s_c . For the decoupling method, we slightly abuse the notation $z_{\theta}(s)$ (or simplified as z) to represent the neural network embedded representations of either s_r or s_c in this section.

3.1 MOTIVATION

Manually abstracting representations of reward-related or cost-related aspects directly from orig-199 inal state observations can be challenging due to the entanglement of information within certain 200 observation dimensions. It is observed that the optimal Q-value, whether associated with reward or 201 cost, exclusively encompasses the information it was trained with, independent of another factor. 202 For instance, concerning the states of the agent depicted in Figure 1, the optimal Q-values related 203 to reward are the same across all actions regardless of the cost-related observations. These states should be embedded as the same reward-related representation. To achieve a coarser abstraction 204 while maintaining the optimal Q-value unchanged, we design the objective for both reward- and 205 cost-related representations as follows 206

$$\max_{\theta} H(s|z_{\theta}(s)) \qquad \text{s.t.} \left(\mathcal{B}^* Q(z_{\theta}(s), a) - Q(z_{\theta}(s), a) \right)^2 \le \epsilon_{\mathcal{B}}, \ \forall a \in \mathcal{A},$$
(4)

208 where $H(\cdot|\cdot)$ represents conditional entropy, $\epsilon_{\mathcal{B}}$ is an arbitrary small number and \mathcal{B}^* is the optimal 209 general/safety Bellman operator. We define $d(s_1, s_2) := \sup_{a \in \mathcal{A}} |Q^*(z_{\theta}(s_1), a), Q^*(z_{\theta}(s_2), a)|$ 210 as the distance measure between a pair of states $s_1, s_2 \in S$. One can always find an arbitrarily small number ϵ_d such that the objective in Eq. 4 can be achieved through embedding the states 211 $\mathcal{C}(s') := \{ \tilde{s} \in \mathcal{S} \mid d(\tilde{s}, s') < \epsilon_d \}$ with the same representation for any $s' \in \mathcal{S}'$, where \mathcal{S}' is a 212 smallest subset of S such that for any $s'_1, s'_2 \in S'$, we have $d(s'_1, s'_2) \ge 2\epsilon_d$ and $\bigcup_{s' \in S'} C(s') = S$. 213 Contrastive learning, which aims to bring artificially defined similar instances closer and push other 214 instances further apart in the representation space (Oord et al., 2018; Bachman et al., 2019; Chen 215 et al., 2020; Agarwal et al., 2021), provides a promising solution for our embedding task.

216 217 3.2 Q-SUPERVISED CONTRASTIVE REPRESENTATION

Inspired by Agarwal et al. (2021), we adopt a soft similarity measure, denoted as $\Gamma(s, \tilde{s}) = \exp(-d(s, \tilde{s})/\eta)$, to quantify the distance between two states (with η representing the temperature factor). Notably, directly calculating the distance measure involves querying out-of-distribution (OOD) actions in the offline setting. To address this issue, we pre-train a generative model to capture the behavior policy (cf. Appendix B.2,C.1 for details). This allows us to generate in-support actions for any given states in the offline dataset, denoted as \mathcal{A}^s_{β} . As a result, we have the approximation $d(s, \tilde{s}) \approx \sup_{a \in \mathcal{A}^s_{\alpha}} |Q^*(z_{\theta}(s), a) - Q^*(z_{\theta}(\tilde{s}), a)|$ for calculating the soft similarity measure.

225 In practice, we employ a random sampling approach to select a subset of states, denoted as S', from 226 the offline state set. Within the subset, we further randomly choose a set of anchor states, denoted as 227 $\{s_i \in \mathcal{S}' \mid i \in \mathcal{I}\}$, where \mathcal{I} represents the index set of the selected anchor states. For each anchor 228 state s_i , we use its nearest neighbor in S' based on the similarity measure Γ to define the positive 229 pairs $\{s_i, \tilde{s}_i\}$, where $\tilde{s}_i = \arg \min_{s \in S' \setminus \{s_i\}} \Gamma(s_i, s)$. The remaining states in S' are considered as 230 negative samples. Attention-based or multiple-layer-perceptron-based neural networks are utilized to encode the state as a normalized vector on the unit hypersphere, i.e., $||z_{\theta}(s)|| = 1$ (cf. Appendix C.2 231 for detailed network selection and structure design). Finally, we have the following contrastive loss, 232 which encourages the embedding of states with similar Q^* values across all actions to have similar 233 representations: 234

$$\mathcal{L}_{\theta} = \sum_{i \in \mathcal{I}} -\frac{1}{|\mathcal{I}|} \log \frac{\Gamma(s_i, \tilde{s}_i) \exp(z_i \cdot \tilde{z}_i/\nu)}{\Gamma(s_i, \tilde{s}_i) \exp(z_i \cdot \tilde{z}_i/\nu) + \sum_{z_j \in \mathcal{Z}' \setminus \{z_i, \tilde{z}_i\}} (1 - \Gamma(z_i, z_j)) \exp(z_i \cdot z_j/\nu)}, \quad (5)$$

238 where ν is a temperature parameter.

It is important to note that Eq. 5 requires precise calculation of optimal Q-values for all states across all actions, i.e., the constraints in Eq. 4 are satisfied. However, the Q-values are derived from the representation network, and even a small change in the network can result in variations in the Q-values. Therefore, it is necessary to integrate the training process of the representation network with the Q-learning process. This integration can be achieved by incorporating the contrastive loss as an auxiliary objective during the learning of the optimal value functions. Please refer to Section 3.3 for further details.

246 247

235 236 237

3.3 PRACTICAL IMPLEMENTATION

Building upon in-sample learning methods (Kostrikov et al., 2021; Xu et al., 2023; Garg et al., 2023; Zheng et al., 2024), our approach follows a two-step process. In the initial phase, we undertake the learning process for the value functions and representations associated with cost and reward separately. Following that, we extract the policy based on the acquired value functions and representations.

Reward-related Representation. We use implicit Q-learning (IQL) (Kostrikov et al., 2021) (cf. Appendix B.1 for details) to approximate the reward-related optimal (maximum) value functions Q_r^* and V_r^* within the support of data distribution through expectile regression:

$$\mathcal{L}_{V_r} = \mathbb{E}_{(s,a)\sim D_\beta} \left[L^{\tau}(Q_r(z_{\theta_r}(s), a) - V_r(z_{\theta_r}(s))) \right], \tag{6}$$

256 257 258

263

268

$$\mathcal{L}_{Q_r} = \mathbb{E}_{(s,a,s',r)\sim D_\beta} \left[(r + \gamma V_r(z_{\theta_r}(s')) - Q_r(z_{\theta_r}(s),a))^2 \right],\tag{7}$$

where $L^{\tau}(u) = |\tau - \mathbb{I}(u < 0)|u^2, \tau \in (0.5, 1)$. The reward-related representations are acquired using the auxiliary contrastive loss term described in Eq. 5, with a weighting factor denoted by δ . Consequently, the overall loss for the reward-related value functions and representations is formulated as:

$$\mathcal{L}_{reward} = \mathcal{L}_{V_r} + \mathcal{L}_{Q_r} + \delta \mathcal{L}_{\theta_r}.$$
(8)

Cost-related Representation. Similar to Zheng et al. (2024), we employ the safety Bellman operator, as denoted in Eq. 2, and utilize reverse expectile regression to learn the cost-related optimal (minimum) value functions Q_h^* and V_h^* :

$$\mathcal{L}_{V_h^{\text{low}}} = \mathbb{E}_{(s,a)\sim D_\beta} \left[L_{\text{rev}}^{\tau}(Q_h(z_{\theta_h}(s), a) - V_h^{\text{low}}(z_{\theta_h}(s))) \right],\tag{9}$$

$$\mathcal{L}_{Q_h} = \mathbb{E}_{(s,a,s',h)\sim D_{\beta}} \left[((1-\gamma)h + \gamma \max\{h, V_h^{\text{low}}(z_{\theta_h}(s')\} - Q_h(z_{\theta_h}(s), a))^2 \right], \quad (10)$$

where $L_{rev}^{\tau}(u) = |\tau - \mathbb{I}(u > 0)|u^2, \tau \in (0.5, 1)$. Additionally, we learn an upper-bound cost-related value function V_h^{up} to model the maximum Q_h^* across all actions in support:

$$\mathcal{L}_{V_h^{\mathrm{up}}} = \mathbb{E}_{(s,a)\sim D_\beta} \left[L^\tau(Q_h(z_{\theta_h}(s), a) - V_h^{\mathrm{up}}(z_{\theta_h}(s))) \right].$$
(11)

By incorporating the auxiliary contrastive loss term (Eq. 5) with a weighting factor of δ , we express the overall loss for the cost-related value functions and representations as:

$$\mathcal{L}_{cost} = \mathcal{L}_{V_h^{\text{low}}} + \mathcal{L}_{V_h^{\text{up}}} + \mathcal{L}_{Q_h} + \delta \mathcal{L}_{\theta_h}.$$
 (12)

Policy Extraction. As illustrated in the right subplot of Figure 2, we divide the global policy into three components: the reward policy π_r , which solely depends on the reward-related representation; the cost policy π_h , which solely relies on the cost-related representations; and the tradeoff policy π_{to} , which depends on both. We independently train the three policies using weighted regressed diffusion models, an approach pioneered by Zheng et al. (2024). They present that the optimal policy satisfies $\pi^*(a|z) \propto \pi_\beta(a|z) \cdot w(z, a)$, and the optimal policy can be obtained through weighted training of diffusion models. The weighted loss function for the three policies can be expressed as follows:

$$\begin{cases}
\pi_r : \mathcal{L}_{\pi_r} = \mathbb{E}_{\operatorname{var}} \left[w_r(z_{\theta_r}(s), a) \| \zeta - \zeta_{\psi_r}(a_t, z_{\theta_r}(s), t) \| \right] \\
\pi_h : \mathcal{L}_{\pi_h} = \mathbb{E}_{\operatorname{var}} \left[w_h(z_{\theta_h}(s), a) \| \zeta - \zeta_{\psi_h}(a_t, z_{\theta_h}(s), t) \| \right] \\
\pi_{to} : \mathcal{L}_{\pi_{to}} = \mathbb{E}_{\operatorname{var}} \left[w_{to}(z_{\theta_r}(s), z_{\theta_h}(s), a) \| \zeta - \zeta_{\psi_{to}}(a_t, z_{\theta_r}(s), z_{\theta_h}(s), t) \| \right],
\end{cases}$$
(13)

where var represents the variables involved in the expectation, with $t \sim U(1,T), \zeta \sim \mathcal{N}(0,I)$, and $(s,a) \sim D_{\beta}$. The noised action $a_t = \alpha_t a + \sigma_t \zeta$ satisfies the forward transition distribution $\mathcal{N}(a_t | \alpha_t a, \sigma_t \mathbb{I})$ in the diffusion models, and α_t, σ_t are noised schedules. The weights in Eq. 13 can be denoted as:

295 296

297

273 274

275

276 277 278

279

280

281

282

283

284

285

287

289

290

291

 $\begin{cases} w_{r}(z_{\theta_{r}}(s), a) = \exp(\iota_{r}(Q_{r}(z_{\theta_{r}}(s), a) - V_{r}(z_{\theta_{r}}(s)))) \\ w_{h}(z_{\theta_{h}}(s), a) = \exp(-\iota_{h}(Q_{h}(z_{\theta_{h}}(s), a) - V_{h}(z_{\theta_{h}}(s)))) \\ w_{to}(z_{\theta_{r}}(s), z_{\theta_{h}}(s), a) = \exp(\iota_{to}(Q_{r}(z_{\theta_{r}}(s), a) - V_{r}(z_{\theta_{r}}(s))) \cdot \mathbb{I}_{Q_{h}(z_{\theta_{h}}(s), a) \leq 0}, \end{cases}$ (14)

where ι_r , ι_h and ι_{to} are temperatures that control the behavior regularization strength.

After obtaining $\zeta_{\psi_r} \zeta_{\psi_h}$ and $\zeta_{\psi_{to}}$, the three approximated optimal policies can be sampled through the reverse diffusion chain starting from random Gaussian noise (Ho et al., 2020; Song et al., 2020) (cf. Appendix B.3 for details). During the testing phase, we first perform safety assessments on the cost-related representations. If the assessment verifies absolute safety $(V_h^{\text{low}} \leq V_h^{\text{up}} \leq 0)$, we employ the policy π_r . If the assessment indicates borderline safety $(V_h^{\text{low}} \leq 0 < V_h^{\text{up}})$, we utilize the policy π_{to} . In the case of an unsafe condition $(0 < V_h^{\text{low}} \leq V_h^{\text{up}})$, we rely on the policy π_h . See the right subplot of Figure 2 for details.

305 306

3.4 COMPARISON WITH BISIMULATION

307 Bisimulation has been established as a useful tool for abstracting state representations (Definition 3.1), 308 where states with identical transition and reward/cost functions are grouped together (Givan et al., 2003; Castro & Precup, 2010; Castro, 2020). However, employing bisimulation typically entails 310 an additional step of training a model-based estimator to learn the state transition and reward/cost 311 functions. Notably, the estimation of reward/cost functions becomes particularly challenging when 312 the values are sparsely distributed (Lee et al., 2024). In contrast to such a model-based representation approach, learning the representations based on Q^* (Definition 3.2) eliminates the necessity for 313 estimating the exact model dynamics (Givan et al., 2003; Li et al., 2006). With Θ denoting a generic 314 surjective mapping from the ground-truth state space to representation space, we have the following 315 definitions: 316

Definition 3.1. A bisimulation representation Θ_{bisim} is such that for any action a and any represented state z, $\Theta_{\text{bisim}}(s_1) = \Theta_{\text{bisim}}(s_2)$ implies $r(s_1, a) = r(s_2, a)$ (or $c(s_1, a) = c(s_2, a)$) and $\sum_{s' \in \Theta_{\text{bisim}}^{-1}(z)} P(s'|s_1, a) = \sum_{s' \in \Theta_{\text{bisim}}^{-1}(z)} P(s'|s_2, a)$.

Definition 3.2. A Q^* -irrelevance representation Θ_{Q^*} is such that, $\Theta_{Q^*}(s_1) = \Theta_{Q^*}(s_2)$ implies $Q^*(s_1, a) = Q^*(s_2, a)$ for any action a.

322

For any state representation Θ_1, Θ_2 , we say Θ_1 is finer than Θ_2 , denoted as $\Theta_1 \succeq \Theta_2$, if and only if for any states $s_1, s_2 \in S$, $\Theta_1(s_1) = \Theta_1(s_2)$ implies $\Theta_2(s_1) = \Theta_2(s_2)$. Givan et al. (2003) established the relationship between bisiumulation and the Q^* -irrelevance representations ($\Theta_{\text{bisim}} \succeq \Theta_{Q^*}$) for finite-horizon MDPs with respect to the general Bellman operator. In the following theorem, we extend this relationship to infinite-horizon MDPs and incorporate the safety Bellman operator, as described below.

Theorem 3.1. For any MDP, the optimal Q-functions induced by either the general Bellman operator or the safety Bellman operator satisfy $\Theta_{\text{bisim}} \succeq \Theta_{Q^*}$. The optimal policies derived from both bisimulation representation and Q^* -irrelevance representation are also optimal in the ground MDP, i.e., $\pi^*(\Theta_{\text{bisim}}(s)) = \pi^*(\Theta_{Q^*}(s)) = \pi^*(s)$ for any state $s \in S$.

Theorem 3.1 shows that neither of the representation methods alters the optimal policy, while the bisimulation representation is finer than the Q-based representation, which implies that

$$0 \le H(s|\Theta_{\text{bisim}}(s)) \le H(s|\Theta_{Q^*}(s)).$$
(15)

Refer to Appendix A for a complete proof. Since our primary objective is to maximize the conditioned entropy $H(s|z_{\theta}(s))$, our Q-supervised contrastive learning method theoretically surpasses bisimulation in terms of generalization.

340 4 EXPERIMENT

335

336

337

338

339

342 4.1 EVALUATION ON DSRL BENCHMARK

343 We compare the proposed SDQC with several state-of-the-art baseline safe offline RL algorithms 344 on the DSRL benchmark (Liu et al., 2023a), which provides extensive datasets and environment 345 wrappers for safe offline RL performance evaluation. The baseline algorithms include i) BCO-Lag: 346 A PID-Lagrangian-based method (Stooke et al., 2020) that considers cost threshold based on Batch 347 Constrained Q-learning (BCQ) (Fujimoto et al., 2019), ii) CPQ (Xu et al., 2022a): A constrained 348 Q-updating method that incorporates penalties for OOD actions and unsafe actions, iii) COptiDICE 349 (Lee et al., 2022): A DICE (distribution correction estimation) based Lagrangian method that builds 350 upon offline RL algorithm OptiDICE (Lee et al., 2021), iv) CDT (Liu et al., 2023b): A future cost inference method based on Decision Transformer (DT) (Chen et al., 2021b), v) TREBI (Lin et al., 351 2023): A real-time cost budget inference method on the basis of Diffuser (Janner et al., 2022), vi) 352 FISOR (Zheng et al., 2024): An HJ reachability guided method with diffusion policies that firstly 353 considers the hard constraints in safe offline RL problems. 354

Our ultimate objective is to achieve zero-cost during test, aligning with the framework established by FISOR (Zheng et al., 2024). However, most baseline algorithms struggle to operate effectively under a zero-cost threshold. Consequently, in accordance with FISOR (Zheng et al., 2024), we impose a stringent cost limit of 10 for the Safety-Gymnasium environment and 5 for Bullet-Safety-Gym. We employ the metrics of *normalized return* and *normalized cost* for evaluation, where a normalized cost below 1 signifies a safe operation. The evaluation results[†] are presented in Table 1.

361 The former five baseline algorithms exhibit either significant constraint violations or suboptimal returns when subjected to stringent safety requirements, partially due to the fact that they consider 362 only soft constraints. Despite incorporating hard constraints, FISOR still encounters high costs in 363 tasks with high complexity. As discussed in Section 1, this issue can be attributed to estimation errors 364 and OOD problems during the testing phase. In contrast to FISOR, our proposed SDQC conducts 365 safety assessments on the cost-related representation abstracted from the original observations and 366 makes decisions accordingly. The utilization of decoupled representations in SDQC substantially 367 improves the accuracy of state safety assessment and enhances the generalization capability of 368 the policy, thereby providing a higher level of safety assurance. The experimental results clearly 369 demonstrate that SDQC outperforms FISOR in terms of higher rewards and lower costs. Remarkably, 370 SDQC even achieves zero violations in the majority of tasks.

371 372

4.2 GENERALIZATION TESTS

To showcase the superior generalization capabilities of our proposed SDQC compared to other safe offline RL algorithms, we perform generalization tests on the "CarGoal" and "CarPush" tasks (in

 [†]The baseline algorithm evaluation results are sourced from FISOR (Zheng et al., 2024), except for the
 evaluation of the Point agent on Safety-Gymnasium (marked with *), which is conducted independently as it is not in the source.

380

381

Table 1: Normalized DSRL benchmark results. The evaluation results are averaged over 3 random seeds (20 episodes for each seed). Gray: Unsafe agents. Bold: Safe agents whose normalized cost is smaller than 1. Red: Safe agents with the highest reward. Blue: Safe agents with the lowest cost.

001															
382	Tack	BCQ-	Lag	CP	Q	COptil	DICE	CD	Т	TREBI		FISOR		SDQC(ours)	
000	Idsk	reward↑	cost↓	reward↑	cost↓										
383	*PointGoal1	0.71	4.29	0.56	0.93	0.40	5.53	0.21	1.59	0.36	2.79	0.68	4.19	0.35	0.36
384	*PointGoal2	0.62	3.81	0.41	5.03	0.43	2.78	0.22	1.19	0.28	3.86	0.21	1.42	0.29	0.09
385	*PointPush1	0.32	3.08	0.14	1.35	0.13	3.80	0.27	2.81	0.31	2.02	0.27	1.38	0.12	0.00
	*PointPush2	0.21	1.86	0.16	2.36	0.02	2.90	0.18	1.69	0.13	3.85	0.24	2.41	0.19	0.28
386	*PointButton1	0.21	4.45	0.61	11.80	0.08	4.29	0.48	10.88	0.12	4.55	0.04	0.97	0.08	0.46
387	*PointButton2	0.38	8.04	0.35	12.09	0.17	6.12	0.42	9.97	0.02	2.18	0.08	4.49	0.06	0.57
200	CarGoal1	0.44	2.76	0.33	4.93	0.43	2.81	0.60	3.15	0.41	1.16	0.49	0.83	0.38	0.01
300	CarGoal2	0.34	4.72	0.10	6.31	0.19	2.83	0.45	6.05	0.13	1.16	0.06	0.33	0.23	0.00
389	CarPush1	0.23	1.33	0.08	0.77	0.21	1.28	0.27	2.12	0.26	1.03	0.28	0.28	0.30	0.00
390	CarPush2	0.10	2.78	-0.03	10.00	0.10	4.55	0.16	4.60	0.12	2.65	0.14	0.89	0.31	0.04
	CarButton1	0.13	6.68	0.22	40.06	-0.16	4.63	0.17	7.05	0.07	3.75	-0.02	0.26	0.03	0.32
391	CarButton2	-0.04	4.43	0.08	19.03	-0.17	3.40	0.23	12.87	-0.03	0.97	0.01	0.58	0.02	0.42
392	AntVel	0.85	18.54	-1.01	0.00	1.00	10.29	0.98	0.91	0.31	0.00	0.89	0.00	0.73	0.00
202	HalfCheetah Vel	1.04	57.06	0.08	2.56	0.43	0.00	0.97	0.55	0.87	0.23	0.89	0.00	0.81	0.00
393	SwimmerVel	0.29	4.10	0.31	11.58	0.58	23.64	0.67	1.47	0.42	1.31	-0.04	0.00	-0.04	0.00
394	SafetyGym	0.39	9.17	0.16	9.05	0.26	5.66	0.46	4.93	0.25	2.10	0.28	1.25	0.26	0.17
395	Average	0.67	2.20	0.00	0.00	0.60	2.64	0.70	1.00	0.60	5.10	0.45	0.03	0.01	0.00
	AntRun	0.65	3.30	0.00	0.00	0.62	3.64	0.70	1.88	0.63	5.43	0.45	0.03	0.31	0.00
396	BallRun	0.43	6.25	0.85	13.67	0.55	11.32	0.32	0.45	0.29	4.24	0.18	0.00	0.20	0.00
397	CarRun	0.84	2.51	1.06	10.49	0.92	0.00	0.99	1.10	0.97	1.01	0.73	0.14	0.56	0.00
200	DroneRun	0.80	17.98	0.02	7.95	0.72	13.//	0.58	0.30	0.59	1.41	0.30	0.55	0.30	0.50
390	AntCircle	0.67	19.13	0.00	0.00	0.18	13.41	0.48	7.44	0.37	2.50	0.20	0.00	0.38	0.00
399	BallCircle	0.67	8.50	0.40	4.37	0.70	9.06	0.68	2.10	0.63	1.89	0.34	0.00	0.42	0.00
400	CarCircle	0.68	8.84	0.49	4.48	0.44	1.13	0./1	2.19	0.49	0.73	0.40	0.11	0.50	0.00
	DroneCircle	0.95	18.50	-0.27	1.29	0.24	2.19	0.55	1.29	0.54	2.30	0.48	0.00	0.36	0.07
401	BuiletGym	0.71	10.63	0.32	5.28	0.55	7.64	0.63	2.09	0.56	2.44	0.39	0.10	0.39	0.08
402	Average														

403

404 Safety-Gymnasium), as illustrated in Figure 3. In these tasks, the "Car" agent is tasked with reaching 405 the goal point or pushing the box to the goal point while avoiding hazardous areas and obstacles. 406 The difficulty level varies between tasks, with the simple tasks (CarGoal1, CarPush1) having fewer 407 hazards and obstacles than the challenging tasks (CarGoal2, CarPush2). 408

It is reasonable to be concerned about the performance of an agent when it is tested in environments 409 that differ from the ones it was trained on, especially if the testing environment is more complex or 410 comprehensive. The experimental results provide evidence that our proposed SDQC algorithm is the 411 only algorithm that ensures no increase in cost under such circumstances. In fact, SDQC achieves 412 almost zero violations in the majority of tests, with only a slight decay in reward performance. In 413 contrast, other algorithms exhibit a sharp increase in cost and/or a significant decrease in reward. 414 Generalization in ensuring safety is crucial in safety-critical scenarios, such as the field of autonomous 415 driving. It is impractical to have the agent traverse every possible radar observation that may occur in 416 real-world scenarios during the training process. Our proposed SDQC offers a potent and promising solution for addressing these complex safety-critical scenarios. 417

418 419

420

4.3 ABLATION STUDY

421 To validate the efficacy of our proposed Q-supervised contrastive learning approach in acquiring 422 meaningful representations and enhancing performance, as discussed in Section 3.2, we conducted ablation studies on the "Safety-Gymnasium-CarGoal2" task (cf. Figure 4). In the absence of 423 contrastive loss during the critic and representation training phase, the agent experiences considerably 424 lower rewards and higher costs compared to the agent trained with contrastive loss. 425

426 The t-SNE visualization results (in Figure 4b) reveal that the Q-supervised contrastive loss effectively 427 clusters representations with similar values in high-dimensional space. This aligns with our original 428 intention, which aims to cluster states with similar Q-values for any actions in the representation space. 429 The clusterings facilitate the learning of the conditional diffusion model (i.e., the actor) by promoting the generation of similar output policies for similar representations. Furthermore, the inclusion of the 430 Q-supervised contrastive loss enables a more reliable evaluation of the states' safety. In the depicted 431 10 trajectories, despite the cumulative cost being zero, the agent trained without contrastive loss



Figure 3: The generalization tests on the agent "Car" in Safety-Gymnasium. (a) The agent is trained on the dataset from a simple environment (S-trained), and its performance is evaluated in both the simple environment (S-tested) and the complicated one (C-tested). Conversely, (b) the agent is trained on the dataset from a complicated environment (C-trained), and its performance is assessed in both the original environment (C-tested) and the simple one (S-tested). The evaluation results are obtained from 3 random seeds, with 20 tests on each seed. Outlier data points are omitted for clarity.

erroneously identifies a majority of the experienced states as unsafe $(V_h^{\text{low}} \ge 0)$. Conversely, the agent trained with contrastive loss provides a more accurate assessment, demonstrating the effectiveness of the proposed approach. For more ablation studies on the impact of anchor number $(|\mathcal{I}|)$ and neural network structures, please refer to Appendix C.

461 462 463

464

449

450

451

452

453

454

455 456

457

458

459

460

5 RELATED WORKS

Safe RL. In online settings, safe RL problems are generally tackled with three mainstream ap-465 proaches (Xu et al., 2022b). i) Formulating the problem as a CMDP and solving it from an opti-466 mization perspective. Solution techniques include updating the policy constrained in a trust region 467 (Achiam et al., 2017; Liu et al., 2022), reformulating the problem into its Lagrangian dual form 468 (Tessler et al., 2018; Chow et al., 2018; Ma et al., 2021b; Duan et al., 2022), and addressing the 469 constraints by framing an optimistic/pessimistic planning problem (Wachi et al., 2018; Kalagarla 470 et al., 2021). Due to expectation constraints and estimation errors, optimization-based methods can 471 only achieve soft constraints, leading to possible violations of the cost threshold during the testing phase (Liu et al., 2023a; Zheng et al., 2024). ii) Combining the safe RL problem with the field of 472 safe control. A prevalent method entails representing a safety certificate through a learned safety 473 assessment function, such as the Control Barrier Function (CBF) (Ma et al., 2021a; Luo & Ma, 2021) 474 or Hamilton-Jacobbi (HJ) reachability (Yu et al., 2022; Fisac et al., 2019; Chen et al., 2021a). An 475 agent can switch between optimal and safe policies based on safety assessment results (Chen et al., 476 2021a; Thananjeyan et al., 2021), thereby theoretically ensuring hard constraints with state-wise 477 zero violations. iii) Employing Teacher-Student Framework (TSF), wherein a proficient teacher 478 critic or policy supervises the student policy, offering guidance and intervening during the onset 479 of safety-critical conditions. (Mehta et al., 2020; Peng et al., 2022; Xue et al., 2023). Applying 480 the aforementioned approaches directly in offline settings may give rise to significant distributional 481 shift issues due to the inaccurate estimation of Q-values for OOD states and actions (Liu et al., 482 2023a). Recent research endeavors have embraced the integration of safety constraint problems with existing reliable offline reinforcement learning algorithms (Xu et al., 2022a; Lee et al., 2022; Liu 483 et al., 2023b; Lin et al., 2023). However, most existing methods only provide soft constraints without 484 any guarantees of zero violations. FISOR (Zheng et al., 2024) is the first safe offline RL algorithm 485 that tackles hard constraints issues through HJ reachability analysis, while the limited offline data



Figure 4: Ablation studies on the Q-supervised contrastive loss in CarGoal2. (a) The actor-trainingprocess evaluations of SDQC with (marked by \bigstar) and without (marked by \times) contrastive loss. The curves are averaged over 3 random seeds and smoothed with a window size of 3. (b) t-SNE visualization of the distribution of the original state, the reward-related and the cost-related representations with and without contrastive loss across 10 different safe trajectories, where the policies are from the agent trained with contrastive loss. The original states (first column) are colored according to the critic trained with contrastive loss (i.e., the same as the second column).

still makes it difficult to guarantee safety during tests thoroughly. As a complementary algorithm to FISOR, our SDQC decouples the global observations for safe decision-making, substantially improves the accuracy of state safety assessment, and enhances the generalization capability of the policy, thereby providing a higher level of safety assurance.

- 516 **Representation Learning.** Representation learning in RL involves compressing the large observa-517 tion space into a smaller latent vector that captures relevant aspects of the environment (Watter et al., 2015; Finn et al., 2016; Gelada et al., 2019), often applied in image-based tasks (Kostrikov et al., 518 2020; Yarats et al., 2021; Cetin et al., 2022). Contrastive learning has been widely acknowledged 519 as a potent technique for unsupervised representation learning (Liu et al., 2021; Zhu et al., 2022), 520 primarily achieved by augmenting data through introducing noise to the original image (Laskin 521 et al., 2020; Agarwal et al., 2021). In state-based tasks, this approach is not directly applicable as 522 the noise may distort the underlying information. Unlike previous works that conduct contrastive 523 learning among the generated samples, we employ contrastive learning within the dataset itself in a 524 Q-supervised manner. The most relevant works to ours are from Bellemare et al. (2019) and Le Lan 525 et al. (2021), who learn representations via Bellman value functions. To the best of our knowledge, 526 we are pioneers in utilizing representation learning in state-based Safe RL tasks. We are the first to 527 introduce the concept of decoupling states into reward- and cost-related representations specifically 528 for decision-making purposes.
- 529 530

6 CONCLUSION

531 532

In this work, we propose the first framework of state decoupling for safe decision-making to tackle the OOD problem of offline safe RL during the testing phase. We propose a Q-supervised contrastive learning method to learn the representations without relying on additional system model estimation such as bisimulation. Theoretical analysis demonstrates that our Q-supervised approach generates coarser representations while preserving the optimal policy, leading to enhanced generalization performance. Experiments on DSRL benchmarks showcase that SDQC surpasses other baseline algorithms, especially for its exceptional ability to achieve almost zero violations in more than half of tasks. Further, SDQC possesses superior generalization ability when confronted with unseen, even more complex environments.

540 REFERENCES

- Joshua Achiam, David Held, Aviv Tamar, and Pieter Abbeel. Constrained policy optimization. In
 International Conference on Machine Learning, pp. 22–31. PMLR, 2017.
- Rishabh Agarwal, Marlos C Machado, Pablo Samuel Castro, and Marc G Bellemare. Contrastive
 behavioral similarity embeddings for generalization in reinforcement learning. *arXiv preprint arXiv:2101.05265*, 2021.
- Eitan Altman. Constrained markov decision processes with total cost criteria: Lagrangian approach and dual linear program. *Mathematical Methods of Operations Research*, 48:387–417, 1998.
- 550 Eitan Altman. Constrained Markov decision processes. Routledge, 2021.
- Per-Arne Andersen, Morten Goodwin, and Ole-Christoffer Granmo. Towards safe reinforcementlearning in industrial grid-warehousing. *Information Sciences*, 537:467–484, 2020.
- Philip Bachman, R Devon Hjelm, and William Buchwalter. Learning representations by maximizing
 mutual information across views. *Advances in Neural Information Processing Systems*, 32, 2019.
- Marc Bellemare, Will Dabney, Robert Dadashi, Adrien Ali Taiga, Pablo Samuel Castro, Nicolas
 Le Roux, Dale Schuurmans, Tor Lattimore, and Clare Lyle. A geometric perspective on optimal
 representations for reinforcement learning. *Advances in Neural Information Processing Systems*, 32, 2019.
- Lukas Brunke, Melissa Greeff, Adam W Hall, Zhaocong Yuan, Siqi Zhou, Jacopo Panerati, and
 Angela P Schoellig. Safe learning in robotics: From learning-based control to safe reinforcement
 learning. Annual Review of Control, Robotics, and Autonomous Systems, 5:411–444, 2022.
- Pablo Castro and Doina Precup. Using bisimulation for policy transfer in MDPs. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 24, pp. 1065–1070, 2010.
- Pablo Samuel Castro. Scalable methods for computing state similarity in deterministic Markov decision processes. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 34, pp. 10069–10076, 2020.
- Edoardo Cetin, Philip J Ball, Steve Roberts, and Oya Celiktutan. Stabilizing off-policy deep
 reinforcement learning from pixels. *arXiv preprint arXiv:2207.00986*, 2022.
- ⁵⁷²
 ⁵⁷³ Bingqing Chen, Jonathan Francis, Jean Oh, Eric Nyberg, and Sylvia L Herbert. Safe autonomous racing via approximate reachability on ego-vision. *arXiv preprint arXiv:2110.07699*, 2021a.
- Lili Chen, Kevin Lu, Aravind Rajeswaran, Kimin Lee, Aditya Grover, Misha Laskin, Pieter Abbeel,
 Aravind Srinivas, and Igor Mordatch. Decision transformer: Reinforcement learning via sequence
 modeling. Advances in Neural Information Processing Systems, 34:15084–15097, 2021b.
- Ting Chen, Simon Kornblith, Mohammad Norouzi, and Geoffrey Hinton. A simple framework for contrastive learning of visual representations. In *International Conference on Machine Learning*, pp. 1597–1607. PMLR, 2020.
- 581
 582 Yinlam Chow, Mohammad Ghavamzadeh, Lucas Janson, and Marco Pavone. Risk-constrained reinforcement learning with percentile risk criteria. *Journal of Machine Learning Research*, 18 (167):1–51, 2018.
- Jingliang Duan, Zhengyu Liu, Shengbo Eben Li, Qi Sun, Zhenzhong Jia, and Bo Cheng. Adaptive
 dynamic programming for nonaffine nonlinear optimal control problem with state constraints.
 Neurocomputing, 484:128–141, 2022.
- Chelsea Finn, Xin Yu Tan, Yan Duan, Trevor Darrell, Sergey Levine, and Pieter Abbeel. Deep spatial autoencoders for visuomotor learning. In 2016 IEEE International Conference on Robotics and Automation (ICRA), pp. 512–519. IEEE, 2016.
- Jaime F Fisac, Neil F Lugovoy, Vicenç Rubies-Royo, Shromona Ghosh, and Claire J Tomlin. Bridging
 Hamilton-Jacobi safety analysis and reinforcement learning. In 2019 International Conference on Robotics and Automation (ICRA), pp. 8550–8556. IEEE, 2019.

 Javier Garcia and Fernando Fernández. A comprehensive survey on safe reinforcement learning. <i>Journal of Machine Learning, Research</i>, 16(1):1437–1480, 2015. Divyansh Garg, Joey Hejna, Matthieu Geist, and Stefano Ermon. Extreme Q-learning: Maxent RL without entropy. <i>arXiv preprint arXiv:2301.02328</i>, 2023. Carles Gelada, Saurabh Kumar, Jacob Buckman, Oir Nachum, and Marc G Bellemare. DeepMDP: Learning continuous latent space models for representation learning. In <i>International Conference on Machine Learning</i>, pp. 2170–2179. PMLR, 2019. Robert Givan, Thomas Dean, and Matthew Greig. Equivalence notions and model minimization in markov decision processes. <i>Artificial Intelligence</i>, 147(1-2):163–223, 2003. Sven Gronauer. Bullet-safety-gym: A framework for constrained reinforcement learning. 2022. Shangding Gu, Long Yang, Yali Du, Guang Chen, Florian Walter, Jun Wang, Yaodong Yang, and Alois Knoll. A review of safe reinforcement learning: Methods, theory and applications. <i>arXiv preprint arXiv:2205.10330</i>, 2022. Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. <i>Advances in Neural Information Processing Systems</i>, 33:6840–6851, 2020. Michael Janner, Yilun Du, Joshua B Tenenbaum, and Sergey Levine. Planning with diffusion for flexible behavior synthesis. <i>arXiv preprint arXiv:2205.09991</i>, 2022. Krishna C Kalagarla, Rahul Jain, and Pietrhigi Nuzzo. A sample-efficient algorithm for episodic finitehorizon mdp with constraints. In <i>Proceedings of the AAAI Conference on Artificial Intelligence</i>, volume 35, pp. 8030–8037, 2021. Jaya Kostrikov, Oenis Yarats, and Rob Fergus. Image augmentation is all you need: Regularizing deep reinforcement learning from pixels. <i>arXiv preprint arXiv:2204.13649</i>, 2020. Jya Kostrikov, Oenis Yarats, and Rob Fergus. Image augmentation is all you need: Regularizing deep reinforcement learning. In <i>International Conferenc</i>	594 595	Scott Fujimoto, David Meger, and Doina Precup. Off-policy deep reinforcement learning without exploration. In <i>International Conference on Machine Learning</i> , pp. 2052–2062. PMLR, 2019.
 Divyansh Garg, Joey Hejna, Mathieu Geist, and Stefano Ermon. Extreme Q-learning: Maxent RL without entropy. <i>arXiv preprint arXiv:2301.02328</i>, 2023. Carles Gelada, Saurabh Kumar, Jacob Buckman, Ofir Nachum, and Marc G Bellemare. DeepMDP: Learning continuous latent space models for representation learning. In <i>International Conference on Machine Learning</i>, pp. 2170–2179. PMLR, 2019. Robert Givan, Thomas Dean, and Matthew Greig. Equivalence notions and model minimization in markov decision processes. <i>Artificial Intelligence</i>, 147(1-2):163–223, 2003. Sven Gronauer. Bullet-safety-gym: A framework for constrained reinforcement learning. 2022. Shangding Gu, Long Yang, Yali Du, Guang Chen, Florian Walter, Jun Wang, Yaodong Yang, and Alois Knoll. A review of safe reinforcement learning: Methods, theory and applications. <i>arXiv preprint arXiv:2205.10330</i>, 2022. Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. <i>Advances in Neural Information Processing Systems</i>, 33:6840–6851, 2020. Michael Janner, Yilun Du, Joshua B Tenenbaum, and Sergey Levine. Planning with diffusion for flexible behavior synthesis. <i>arXiv preprint arXiv:2205.09991</i>, 2022. Krishna C Kalagarla, Rahul Jain, and Pierluigi Nuzzo. A sample-efficient algorithm for episodic finite-horizon mdp with constraints. In <i>Proceedings of the AAAI Conference on Artificial Intelligence</i>, volume 35, pp. 8030–8037, 2021. B Ravi Kiran, Ibrahim Sobh, Victor Talpaert, Patrick Mannion, Ahmad A Al Sallab, Senthi Yogamani, and Patrick Pérez. Deep reinforcement learning for autonomous driving: A survey. <i>IEEE Transactions on Intelligent Transportation Systems</i>, 23(6):4909–4926, 2021. Ilya Kostrikov, Ashvin Nair, and Sergey Luvine. Offline reinforcement learning with implicit Q-learning, <i>arXiv preprint arXiv:2110.06169</i>, 2021. Aviral Kumar, Justin Fu, Matthew Soh, George Tucker, and Sergey Levine. Stab	596 597 598	Javier Garcia and Fernando Fernández. A comprehensive survey on safe reinforcement learning. <i>Journal of Machine Learning Research</i> , 16(1):1437–1480, 2015.
 Divyansh Garg, Joey Hejna, Mattheu Gest, and Stefano Ermon. Extreme Q-learning: Makent RL without entropy. <i>arXiv preprint arXiv:2301.02328</i>, 2023. Carles Gelada, Saurabh Kumar, Jacob Buckman, Ofir Nachum, and Marc G Bellemare. DeepMDP: Learning continuous latent space models for representation learning. In <i>International Conference on Machine Learning</i>, pp. 2170–2179. PMLR, 2019. Robert Givan, Thomas Dean, and Matthew Greig. Equivalence notions and model minimization in markov decision processes. <i>Artificial Intelligence</i>, 147(1-2):163–223, 2003. Sven Gronauer. Bullet-safety-gym: A framework for constrained reinforcement learning. 2022. Shangding Gu, Long Yang, Yali Du, Guang Chen, Florian Walter, Jun Wang, Yaodong Yang, and Alois Knoll. A review of safe reinforcement learning: Methods, theory and applications. <i>arXiv preprint arXiv:2205.10330</i>, 2022. Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. <i>Advances in Neural Information Processing Systems</i>, 33:6840–6851, 2020. Michael Janner, Yilun Du, Joshua B Tenenbaum, and Sergey Levine. Planning with diffusion for flexible behavior synthesis. <i>arXiv preprint arXiv:2205.09991</i>, 2022. Krishna C Kalagarla, Rahul Jain, and Pietur Japter, Patrick Mannion, Ahmad A Al Sallab, Senthil Yogamani, and Patrick Pérez. Deep reinforcement learning for autonomous driving: A survey. <i>IEEE Transactions on Intelligent Transportation Systems</i>, 23(6):4909–4926, 2021. Ilya Kostrikov, Denis Yarats, and Rob Fergus. Image augmentation is all you need: Regularizing deep reinforcement learning from pixels. <i>arXiv preprint arXiv:204.13649</i>, 2020. Hya Kostrikov, Dehis Yarats, and Rob Fergus. Image augmentation is all you need: Regularizing deep reinforcement learning from pixels. <i>arXiv preprint arXiv:204.13649</i>, 2020. Hya Kostrikov, Dehis Yarats, and Rob Fergus. Image augmentation is all you need: Regularizing dee	599	
 Carles Gelada, Saurabh Kumar, Jacob Buckman, Ofir Nachum, and Marc G Bellemare. DeepMDP: Learning continuous latent space models for representation learning. In <i>International Conference</i> <i>on Machine Learning</i>, pp. 2170–2179. PMLR, 2019. Robert Givan, Thomas Dean, and Matthew Greig. Equivalence notions and model minimization in markov decision processes. <i>Artificial Intelligence</i>, 147(1-2):163–223, 2003. Sven Gronauer. Bullet-safety-gym: A framework for constrained reinforcement learning. 2022. Shangding Gu, Long Yang, Yali Du, Guang Chen, Florian Walter, Jun Wang, Yaodong Yang, and Alois Knoll. A review of safe reinforcement learning: Methods, theory and applications. <i>arXiv</i> <i>preprint arXiv:2205.10330</i>, 2022. Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. <i>Advances in</i> <i>Neural Information Processing Systems</i>, 33:6840–6851, 2020. Michael Janner, Yilun Du, Joshua B Tenenbaum, and Sergey Levine. Planning with diffusion for flexible behavior synthesis. <i>arXiv preprint arXiv:2205.09991</i>, 2022. Krishna C Kalagarla, Rahul Jain, and Pierluigi Nuzzo. A sample-efficient algorithm for episodic finite- horizon mdp with constraints. In <i>Proceedings of the AAAI Conference on Artificial Intelligence</i>, volume 35, pp. 8030–8037, 2021. B Ravi Kiran, Ibrahim Sobh, Victor Talpaert, Patrick Mannion, Ahmad A Al Sallab, Senthil Yogamani, and Patrick Pérez. Deep reinforcement learning for autonomous driving: A survey. <i>IEEE</i> <i>Transactions on Intelligent Transportation Systems</i>, 23(6):4909–4926, 2021. Ilya Kostrikov, Denis Yaats, and Rob Fergus. Image augmentation is all you need: Regularizing deep reinforcement learning from pixels. <i>arXiv preprint arXiv:2204.13649</i>, 2020. Ilya Kostrikov, Ashvin Nair, and Sergey Levine. Offline reinforcement learning with implicit Q-learning via bootstrapping error reduction. <i>Advances in Neural Information Processing Systems</i>, 32, 20	600	Divyansh Garg, Joey Hejna, Matthieu Geist, and Stefano Ermon. Extreme Q-learning: Maxent RL without entropy. <i>arXiv preprint arXiv:2301.02328</i> , 2023.
 Carles Getada, Sadiadous Jatenta, Sade Duckman, Oni Vachani, and Mate O Defendate. DeepMD1. Learning continuous latent space models for representation learning. In <i>International Conference on Machine Learning</i>, pp. 2170–2179. PMLR, 2019. Robert Givan, Thomas Dean, and Matthew Greig. Equivalence notions and model minimization in markov decision processes. <i>Artificial Intelligence</i>, 147(1-2):163–223, 2003. Sven Gronauer. Bullet-safety-gym: A framework for constrained reinforcement learning. 2022. Shangding Gu, Long Yang, Yali Du, Guang Chen, Florian Walter, Jun Wang, Yaodong Yang, and Alois Knoll. A review of safe reinforcement learning: Methods, theory and applications. <i>arXiv preprint arXiv:2205.10330</i>, 2022. Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. <i>Advances in Neural Information Processing Systems</i>, 33:6840–6851, 2020. Michael Janner, Yilun Du, Joshua B Tenenbaum, and Sergey Levine. Planning with diffusion for flexible behavior synthesis. <i>arXiv preprint arXiv:2205.09991</i>, 2022. Krishna C Kalagarla, Rahul Jain, and Pierluigi Nuzzo. A sample-efficient algorithm for episodic finite-horizon mdp with constraints. In <i>Proceedings of the AAAI Conference on Artificial Intelligence</i>, volume 35, pp. 8030–8037, 2021. B Ravi Kiran, Ibrahim Sobh, Victor Talpaert, Patrick Mannion, Ahmad A AI Sallab, Senthil Yogamani, and Patrick Pérez. Deep reinforcement learning for autonomous driving: A survey. <i>IEEE Transactions on Intelligent Transportation Systems</i>, 23(6):4909–4926, 2021. Hya Kostrikov, Denis Yarats, and Rob Fergus. Image augmentation is all you need: Regularizing deep reinforcement learning from pixels. <i>arXiv preprint arXiv:2004.13649</i>, 2020. Hya Kostrikov, Ashvin Nair, and Sergey Levine. Offline reinforcement learning with implicit Q-learning via bootstrapping error reduction. <i>Advances in Neural Information Processing Systems</i>, 32, 2019.<td>601</td><td>Corles Calada, Sourabh Kumar, Jacob Buckman, Ofir Nachum, and Mara C. Ballamara, DeenMDP:</td>	601	Corles Calada, Sourabh Kumar, Jacob Buckman, Ofir Nachum, and Mara C. Ballamara, DeenMDP:
 ob Mathie Edminis, pp. 110 Status 2013. Robert Givan, Thomas Dean, and Matthew Greig. Equivalence notions and model minimization in markov decision processes. <i>Artificial Intelligence</i>, 147(1-2):163–223, 2003. Sven Gronauer. Bullet-safety-gym: A framework for constrained reinforcement learning. 2022. Shangding Gu, Long Yang, Yali Du, Guang Chen, Florian Walter, Jun Wang, Yaodong Yang, and Alois Knoll. A review of safe reinforcement learning: Methods, theory and applications. <i>arXiv preprint arXiv:2205.10330</i>, 2022. Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. <i>Advances in Neural Information Processing Systems</i>, 33:6840–6851, 2020. Michael Janner, Yilun Du, Joshua B Tenenbaum, and Sergey Levine. Planning with diffusion for flexible behavior synthesis. <i>arXiv preprint arXiv:2205.09991</i>, 2022. Krishna C Kalagarla, Rahul Jain, and Pietruigi Nuzzo. A sample-efficient algorithm for episodic finite-horizon mdp with constraints. In <i>Proceedings of the AAAI Conference on Artificial Intelligence</i>, volume 35, pp. 8030–8037, 2021. B Ravi Kiran, Ibrahim Sobh, Victor Talpaert, Patrick Mannion, Ahmad A Al Sallab, Senthil Yogamani, and Patrick Pérez. Deep reinforcement learning for autonomous driving: A survey. <i>IEEE Transactions on Intelligent Transportation Systems</i>, 23(6):4909–4926, 2021. Ilya Kostrikov, Denis Yarats, and Rob Fergus. Image augmentation is all you need: Regularizing deep reinforcement learning from pixels. <i>arXiv preprint arXiv:2004.13649</i>, 2020. Michael Laskin, Aravind Srinivas, and Pieter Abbeel. CURL: Contrastive unsupervised representations for reinforcement learning. In <i>International Conference on Artificial Intelligence</i>, volume 35, pp. 8261–8269, 2021. Aviral Kumar, Justin Fu, Matthew Soh, George Tucker, and Sergey Levine. Stabilizing off-policy q-learning. <i>arXiv preprint arXiv:2100.06169</i>, 2021. Michael La	602 603	Learning continuous latent space models for representation learning. In International Conference on Machine Learning, pp. 2170–2179, PMLR, 2019
 Robert Givan, Thomas Dean, and Matthew Greig. Equivalence notions and model minimization in markov decision processes. <i>Artificial Intelligence</i>, 147(1-2):163–223, 2003. Sven Gronauer. Bullet-safety-gym: A framework for constrained reinforcement learning. 2022. Shangding Gu, Long Yang, Yali Du, Guang Chen, Florian Walter, Jun Wang, Yaodong Yang, and Alois Knoll. A review of safe reinforcement learning: Methods, theory and applications. <i>arXiv preprint arXiv:2205.10330</i>, 2022. Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. <i>Advances in Neural Information Processing Systems</i>, 33:6840–6851, 2020. Michael Janner, Yilun Du, Joshua B Tenenbaum, and Sergey Levine. Planning with diffusion for flexible behavior synthesis. <i>arXiv preprint arXiv:2205.09991</i>, 2022. Krishna C Kalagarla, Rahul Jain, and Pietruigi Nuzzo. A sample-efficient algorithm for episodic finite-horizon mdp with constraints. In <i>Proceedings of the AAAI Conference on Artificial Intelligence</i>, volume 35, pp. 8030–8037, 2021. B Ravi Kiran, Ibrahim Sobh, Victor Talpaert, Patrick Mannion, Ahmad A Al Sallab, Senthil Yogamani, and Patrick Pérez. Deep reinforcement learning for autonomous driving: A survey. <i>IEEE Transactions on Intelligent Transportation Systems</i>, 23(6):4909–4926, 2021. Ilya Kostrikov, Ashvin Nair, and Sergey Levine. Offline reinforcement learning with implicit Q-learning. <i>arXiv preprint arXiv:2110.06169</i>, 2021. Aviral Kumar, Justin Fu, Matthew Soh, George Tucker, and Sergey Levine. Stabilizing off-policy q-learning via bootstrapping error reduction. <i>Advances in Neural Information Processing Systems</i>, 32, 2019. Michael Lashin, Aravind Srinivas, and Pieter Abbeel. CURL: Contrastive unsupervised representations for reinforcement learning. In <i>International Conference on Artificial Intelligence</i>, volume 35, pp. 8261–8269, 2021. Jongmin Lee, Osnisok Jeon, Byungjun Le	604	on machine Learning, pp. 2170-2179. 1 MER, 2019.
 Sven Gronauer. Bullet-safety-gym: A framework for constrained reinforcement learning. 2022. Shangding Gu, Long Yang, Yali Du, Guang Chen, Florian Walter, Jun Wang, Yaodong Yang, and Alois Knoll. A review of safe reinforcement learning: Methods, theory and applications. <i>arXiv preprint arXiv:2205.10330</i>, 2022. Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. <i>Advances in Neural Information Processing Systems</i>, 33:6840–6851, 2020. Michael Janner, Yilun Du, Joshua B Tenenbaum, and Sergey Levine. Planning with diffusion for flexible behavior synthesis. <i>arXiv preprint arXiv:2205.09991</i>, 2022. Krishna C Kalagarla, Rahul Jain, and Piertuigi Nuzzo. A sample-efficient algorithm for episodic finite-horizon mdp with constraints. In <i>Proceedings of the AAAI Conference on Artificial Intelligence</i>, volume 35, pp. 8030–8037, 2021. B Ravi Kiran, Ibrahim Sobh, Victor Talpaert, Patrick Mannion, Ahmad A Al Sallab, Senthil Yogamani, and Patrick Pérez. Deep reinforcement learning for autonomous driving: A survey. <i>IEEE Transactions on Intelligent Transportation Systems</i>, 23(6):4909–4926, 2021. Ilya Kostrikov, Ashvin Nair, and Sergey Levine. Offline reinforcement learning with implicit Q-learning. <i>arXiv preprint arXiv:2110.06169</i>, 2021. Aviral Kumar, Justin Fu, Matthew Soh, George Tucker, and Sergey Levine. Stabilizing off-policy q-learning via bootstrapping error reduction. <i>Advances in Neural Information Processing Systems</i>, 32, 2019. Michael Laskin, Aravind Srinivas, and Pieter Abbeel. CURL: Contrastive unsupervised representations for reinforcement learning. In <i>International Conference on Artificial Intelligence</i>, volume 35, pp. 8261–8269, 2021. Jongmin Lee, Wonseok Jeon, Byungjun Lee, Joelle Pineau, and Kee-Eung Kim. OptiDICE: Offline policy optimization via stationary distribution correction estimation. In <i>International Conference on Machine Learning</i>, pp. 612	605 606	Robert Givan, Thomas Dean, and Matthew Greig. Equivalence notions and model minimization in markov decision processes. <i>Artificial Intelligence</i> , 147(1-2):163–223, 2003.
 Shangding Gu, Long Yang, Yali Du, Guang Chen, Florian Walter, Jun Wang, Yaodong Yang, and Alois Knoll. A review of safe reinforcement learning: Methods, theory and applications. <i>arXiv</i> <i>preprint arXiv:2205.10330</i>, 2022. Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. <i>Advances in</i> <i>Neural Information Processing Systems</i>, 33:6840–6851, 2020. Michael Janner, Yilun Du, Joshua B Tenenbaum, and Sergey Levine. Planning with diffusion for flexible behavior synthesis. <i>arXiv preprint arXiv:2205.00991</i>, 2022. Krishna C Kalagarla, Rahul Jain, and Pierluigi Nuzzo. A sample-efficient algorithm for episodic finite- horizon mdp with constraints. In <i>Proceedings of the AAAI Conference on Artificial Intelligence</i>, volume 35, pp. 8030–8037, 2021. B Ravi Kiran, Ibrahim Sobh, Victor Talpaert, Patrick Mannion, Ahmad A Al Sallab, Senthil Yogamani, and Patrick Pérez. Deep reinforcement learning for autonomous driving: A survey. <i>IEEE</i> <i>Transactions on Intelligent Transportation Systems</i>, 23(6):4909–4926, 2021. Ilya Kostrikov, Denis Yarats, and Rob Fergus. Image augmentation is all you need: Regularizing deep reinforcement learning from pixels. <i>arXiv preprint arXiv:2004.13649</i>, 2020. Ilya Kostrikov, Ashvin Nair, and Sergey Levine. Offline reinforcement learning with implicit Q-learning. <i>arXiv preprint arXiv:2110.06169</i>, 2021. Aviral Kumar, Justin Fu, Matthew Soh, George Tucker, and Sergey Levine. Stabilizing off-policy q-learning to bootstrapping error reduction. <i>Advances in Neural Information Processing Systems</i>, 32, 2019. Michael Laskin, Aravind Srinivas, and Pieter Abbeel. CURL: Contrastive unsupervised representa- tions for reinforcement learning. In <i>International Conference on Machine Learning</i>, pp. 5639–5650. PMLR, 2020. Charline Le Lan, Marc G Bellemare, and Pablo Samuel Castro. Metrics and continuity in reinforce- ment learning. In <i>Proceedings of the AAAI </i>	607 608	Sven Gronauer. Bullet-safety-gym: A framework for constrained reinforcement learning. 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. Krishna C Kalagarla, Rahul Jain, and Pierluigi Nuzzo. A sample-efficient algorithm for episodic finite- horizon mdp with constraints. In Proceedings of the AAAI Conference on Artificial Intelligence, volume 35, pp. 8030–8037, 2021. B Ravi Kiran, Ibrahim Sobh, Victor Talpaert, Patrick Mannion, Ahmad A Al Sallab, Senthil Yogamani, and Patrick Pérez. Deep reinforcement learning for autonomous driving: A survey. IEEE Transactions on Intelligent Transportation Systems, 23(6):4909–4926, 2021. Ilya Kostrikov, Denis Yarats, and Rob Fergus. Image augmentation is all you need: Regularizing deep reinforcement learning from pixels. arXiv preprint arXiv:2004.13649, 2020. 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. Michael Laskin, Aravind Srinivas, and Pieter Abbeel. CURL: Contrastive unsupervised representa- tions for reinforcement learning. In International Conference on Machine Learning, pp. 5639–5650. PMLR, 2020. Charline Le Lan, Marc G Bellemare, and Pablo Samuel Castro. Metrics and continuity in reinforce- ment learning. In Proceedings of the AAAI Conference on Artificial Intelligence, volume 35, pp. 8261–8269, 2021. Jongmin Lee, Wonseok Jeon, Byungjun Lee, Joelle Pineau, and Kee-Eung Kim. OptiDICE: Offline policy optimization via stationary distribution correction estimation. In Inte	609 610 611	Shangding Gu, Long Yang, Yali Du, Guang Chen, Florian Walter, Jun Wang, Yaodong Yang, and Alois Knoll. A review of safe reinforcement learning: Methods, theory and applications. <i>arXiv</i> preprint arXiv:2205.10330, 2022.
 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. Krishna C Kalagarla, Rahul Jain, and Pierluigi Nuzzo. A sample-efficient algorithm for episodic finite- horizon mdp with constraints. In <i>Proceedings of the AAAI Conference on Artificial Intelligence</i>, volume 35, pp. 8030–8037, 2021. B Ravi Kiran, Ibrahim Sobh, Victor Talpaert, Patrick Mannion, Ahmad A AI Sallab, Senthil Yogamani, and Patrick Pérez. Deep reinforcement learning for autonomous driving: A survey. <i>IEEE Transactions on Intelligent Transportation Systems</i>, 23(6):4909–4926, 2021. Ilya Kostrikov, Denis Yarats, and Rob Fergus. Image augmentation is all you need: Regularizing deep reinforcement learning from pixels. arXiv preprint arXiv:2004.13649, 2020. Ilya Kostrikov, Ashvin Nair, and Sergey Levine. Offline reinforcement learning with implicit Q-learning via bootstrapping error reduction. Advances in Neural Information Processing Systems, 32, 2019. Michael Laskin, Aravind Srinivas, and Pieter Abbeel. CURL: Contrastive unsupervised representa- tions for reinforcement learning. In International Conference on Machine Learning, pp. 5639–5650. PMLR, 2020. Charline Le Lan, Marc G Bellemare, and Pablo Samuel Castro. Metrics and continuity in reinforce- ment learning. In Proceedings of the AAAI Conference on Artificial Intelligence, volume 35, pp. 8261–8269, 2021. Jongmin Lee, Wonseok Jeon, Byungjun Lee, Joelle Pineau, and Kee-Eung Kim. OptiDICE: Offline policy optimization via stationary distribution correction estimation. In International Conference on Machine Learning, pp. 6120–6130. PMLR, 2021. Jongmin Lee, Cosmin Paduraru, Daniel J Mankowitz, Nicolas Heess, Doina Precup, Kee-Eung Kim, and Arthur Guez. COptiDICE: Offline constrained reinforcement learning via st	612	Jonathan Ho, Aiay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. Advances in
 Michael Janner, Yilun Du, Joshua B Tenenbaum, and Sergey Levine. Planning with diffusion for flexible behavior synthesis. <i>arXiv preprint arXiv:2205.09991</i>, 2022. Krishna C Kalagarla, Rahul Jain, and Pierluigi Nuzzo. A sample-efficient algorithm for episodic finite- horizon mdp with constraints. In <i>Proceedings of the AAAI Conference on Artificial Intelligence</i>, volume 35, pp. 8030–8037, 2021. B Ravi Kiran, Ibrahim Sobh, Victor Talpaert, Patrick Mannion, Ahmad A Al Sallab, Senthil Yogamani, and Patrick Pérez. Deep reinforcement learning for autonomous driving: A survey. <i>IEEE Transactions on Intelligent Transportation Systems</i>, 23(6):4909–4926, 2021. Ilya Kostrikov, Denis Yarats, and Rob Fergus. Image augmentation is all you need: Regularizing deep reinforcement learning from pixels. <i>arXiv preprint arXiv:2004.13649</i>, 2020. Ilya Kostrikov, Ashvin Nair, and Sergey Levine. Offline reinforcement learning with implicit Q-learning. <i>arXiv preprint arXiv:2110.06169</i>, 2021. Aviral Kumar, Justin Fu, Matthew Soh, George Tucker, and Sergey Levine. Stabilizing off-policy q-learning via bootstrapping error reduction. <i>Advances in Neural Information Processing Systems</i>, 32, 2019. Michael Laskin, Aravind Srinivas, and Pieter Abbeel. CURL: Contrastive unsupervised representa- tions for reinforcement learning. In <i>International Conference on Machine Learning</i>, pp. 5639–5650. PMLR, 2020. Charline Le Lan, Marc G Bellemare, and Pablo Samuel Castro. Metrics and continuity in reinforce- ment learning. In <i>Proceedings of the AAAI Conference on Artificial Intelligence</i>, volume 35, pp. 8261–8269, 2021. Jongmin Lee, Wonseok Jeon, Byungjun Lee, Joelle Pineau, and Kee-Eung Kim. OptiDICE: Offline policy optimization via stationary distribution correction estimation. In <i>International Conference on Machine Learning</i>, pp. 6120–6130. PMLR, 2021. Jongmin Lee, Cosmin Paduraru, Daniel J Mankowitz, Nicolas Heess, Doin	613 614	Neural Information Processing Systems, 33:6840–6851, 2020.
 flexible behavior synthesis. <i>arXiv preprint arXiv:2205.09991</i>, 2022. Krishna C Kalagarla, Rahul Jain, and Pierluigi Nuzzo. A sample-efficient algorithm for episodic finite- horizon mdp with constraints. In <i>Proceedings of the AAAI Conference on Artificial Intelligence</i>, volume 35, pp. 8030–8037, 2021. B Ravi Kiran, Ibrahim Sobh, Victor Talpaert, Patrick Mannion, Ahmad A Al Sallab, Senthil Yogamani, and Patrick Pérez. Deep reinforcement learning for autonomous driving: A survey. <i>IEEE Transactions on Intelligent Transportation Systems</i>, 23(6):4909–4926, 2021. Ilya Kostrikov, Denis Yarats, and Rob Fergus. Image augmentation is all you need: Regularizing deep reinforcement learning from pixels. <i>arXiv preprint arXiv:2004.13649</i>, 2020. Ilya Kostrikov, Ashvin Nair, and Sergey Levine. Offline reinforcement learning with implicit Q-learning. <i>arXiv preprint arXiv:2110.06169</i>, 2021. Aviral Kumar, Justin Fu, Matthew Soh, George Tucker, and Sergey Levine. Stabilizing off-policy q-learning via bootstrapping error reduction. <i>Advances in Neural Information Processing Systems</i>, 32, 2019. Michael Laskin, Aravind Srinivas, and Pieter Abbeel. CURL: Contrastive unsupervised representa- tions for reinforcement learning. In <i>International Conference on Machine Learning</i>, pp. 5639–5650. PMLR, 2020. Charline Le Lan, Marc G Bellemare, and Pablo Samuel Castro. Metrics and continuity in reinforce- ment learning. In <i>Proceedings of the AAAI Conference on Artificial Intelligence</i>, volume 35, pp. 8261–8269, 2021. Jongmin Lee, Wonseok Jeon, Byungjun Lee, Joelle Pineau, and Kee-Eung Kim. OptiDICE: Offline policy optimization via stationary distribution correction estimation. In <i>International Conference on Machine Learning</i>, pp. 6120–6130. PMLR, 2021. Jongmin Lee, Cosmin Paduraru, Daniel J Mankowitz, Nicolas Heess, Doina Precup, Kee-Eung Kim, and Arthur Guez. COptiDICE: Offline constrained reinforcement learnin	615	Michael Janner, Yilun Du, Joshua B Tenenbaum, and Sergey Levine. Planning with diffusion for
 Krishna C Kalagarla, Rahul Jain, and Pierluigi Nuzzo. A sample-efficient algorithm for episodic finite- horizon mdp with constraints. In <i>Proceedings of the AAAI Conference on Artificial Intelligence</i>, volume 35, pp. 8030–8037, 2021. B Ravi Kiran, Ibrahim Sobh, Victor Talpaert, Patrick Mannion, Ahmad A Al Sallab, Senthil Yogamani, and Patrick Pérez. Deep reinforcement learning for autonomous driving: A survey. <i>IEEE Transactions on Intelligent Transportation Systems</i>, 23(6):4909–4926, 2021. Ilya Kostrikov, Denis Yarats, and Rob Fergus. Image augmentation is all you need: Regularizing deep reinforcement learning from pixels. <i>arXiv preprint arXiv:2004.13649</i>, 2020. Ilya Kostrikov, Ashvin Nair, and Sergey Levine. Offline reinforcement learning with implicit Q-learning. <i>arXiv preprint arXiv:2110.06169</i>, 2021. Aviral Kumar, Justin Fu, Matthew Soh, George Tucker, and Sergey Levine. Stabilizing off-policy q-learning via bootstrapping error reduction. <i>Advances in Neural Information Processing Systems</i>, 32, 2019. Michael Laskin, Aravind Srinivas, and Pieter Abbeel. CURL: Contrastive unsupervised representa- tions for reinforcement learning. In <i>International Conference on Machine Learning</i>, pp. 5639–5650. PMLR, 2020. Charline Le Lan, Marc G Bellemare, and Pablo Samuel Castro. Metrics and continuity in reinforce- ment learning. In <i>Proceedings of the AAAI Conference on Artificial Intelligence</i>, volume 35, pp. 8261–8269, 2021. Jongmin Lee, Wonseok Jeon, Byungjun Lee, Joelle Pineau, and Kee-Eung Kim. OptiDICE: Offline policy optimization via stationary distribution correction estimation. In <i>International Conference on Machine Learning</i>, pp. 6120–6130. PMLR, 2021. Jongmin Lee, Cosmin Paduraru, Daniel J Mankowitz, Nicolas Heess, Doina Precup, Kee-Eung Kim, and Arthur Guez. COptiDICE: Offline constrained reinforcement learning via stationary distribution correction estimation. <i>arXiv preprint arXiv:2204.089</i>	616	flexible behavior synthesis. arXiv preprint arXiv:2205.09991, 2022.
 Initial of Rangani, Rank Park and Proceedings of the AAAI Conference on Artificial Intelligence, volume 35, pp. 8030–8037, 2021. B Ravi Kiran, Ibrahim Sobh, Victor Talpaert, Patrick Mannion, Ahmad A Al Sallab, Senthil Yogamani, and Patrick Pérez. Deep reinforcement learning for autonomous driving: A survey. <i>IEEE Transactions on Intelligent Transportation Systems</i>, 23(6):4909–4926, 2021. Ilya Kostrikov, Denis Yarats, and Rob Fergus. Image augmentation is all you need: Regularizing deep reinforcement learning from pixels. <i>arXiv preprint arXiv:2004.13649</i>, 2020. Ilya Kostrikov, Ashvin Nair, and Sergey Levine. Offline reinforcement learning with implicit Q-learning. <i>arXiv preprint arXiv:2110.06169</i>, 2021. Aviral Kumar, Justin Fu, Matthew Soh, George Tucker, and Sergey Levine. Stabilizing off-policy q-learning via bootstrapping error reduction. <i>Advances in Neural Information Processing Systems</i>, 32, 2019. Michael Laskin, Aravind Srinivas, and Pieter Abbeel. CURL: Contrastive unsupervised representations for reinforcement learning. In <i>International Conference on Machine Learning</i>, pp. 5639–5650. PMLR, 2020. Charline Le Lan, Marc G Bellemare, and Pablo Samuel Castro. Metrics and continuity in reinforcement learning. In <i>International Conference on Artificial Intelligence</i>, volume 35, pp. 8261–8269, 2021. Jongmin Lee, Wonsek Jeon, Byungjun Lee, Joelle Pineau, and Kee-Eung Kim. OptiDICE: Offline policy optimization via stationary distribution correction estimation. In <i>International Conference on Machine Learning Conference on Machine Learning</i>, pp. 6120–6130. PMLR, 2021. Jongmin Lee, Cosmin Paduraru, Daniel J Mankowitz, Nicolas Heess, Doina Precup, Kee-Eung Kim, and Arthur Guez. COptiDICE: Offline constrained reinforcement learning via stationary distribution correction estimation. <i>arXiv preprint arXiv:2204.08957</i>, 2022. Vint Lee, Pieter Abbeel, and Youngwoon Lee. Dreamsmooth: Improvin	617	Krishna C Kalagarla Rahul Jain and Pierluigi Nuzzo. A sample-efficient algorithm for episodic finite-
 B Ravi Kiran, Ibrahim Sobh, Victor Talpaert, Patrick Mannion, Ahmad A Al Sallab, Senthil Yogamani, and Patrick Pérez. Deep reinforcement learning for autonomous driving: A survey. <i>IEEE</i> <i>Transactions on Intelligent Transportation Systems</i>, 23(6):4909–4926, 2021. Ilya Kostrikov, Denis Yarats, and Rob Fergus. Image augmentation is all you need: Regularizing deep reinforcement learning from pixels. <i>arXiv preprint arXiv:2004.13649</i>, 2020. Ilya Kostrikov, Ashvin Nair, and Sergey Levine. Offline reinforcement learning with implicit Q-learning. <i>arXiv preprint arXiv:2110.06169</i>, 2021. Aviral Kumar, Justin Fu, Matthew Soh, George Tucker, and Sergey Levine. Stabilizing off-policy q-learning via bootstrapping error reduction. <i>Advances in Neural Information Processing Systems</i>, 32, 2019. Michael Laskin, Aravind Srinivas, and Pieter Abbeel. CURL: Contrastive unsupervised representa- tions for reinforcement learning. In <i>International Conference on Machine Learning</i>, pp. 5639–5650. PMLR, 2020. Charline Le Lan, Marc G Bellemare, and Pablo Samuel Castro. Metrics and continuity in reinforce- ment learning. In <i>Proceedings of the AAAI Conference on Artificial Intelligence</i>, volume 35, pp. 8261–8269, 2021. Jongmin Lee, Wonseok Jeon, Byungjun Lee, Joelle Pineau, and Kee-Eung Kim. OptiDICE: Offline policy optimization via stationary distribution correction estimation. In <i>International Conference on Machine Learning</i>, pp. 6120–6130. PMLR, 2021. Jongmin Lee, Cosmin Paduraru, Daniel J Mankowitz, Nicolas Heess, Doina Precup, Kee-Eung Kim, and Arthur Guez. COptiDICE: Offline constrained reinforcement learning via stationary distribution correction estimation. <i>arXiv preprint arXiv:2204.08957</i>, 2022. Vint Lee, Pieter Abbeel, and Youngwoon Lee. Dreamsmooth: Improving model-based reinforcement learning via reward smoothing. In <i>The Twelfth International Conference on Learning Representa-</i> <i>tions 2024</i> 	618 619	horizon mdp with constraints. In <i>Proceedings of the AAAI Conference on Artificial Intelligence</i> , volume 35, pp. 8030–8037, 2021.
 B Ravi Kiran, Ibrahim Sobh, Victor Talpaert, Patrick Mannion, Ahmad A Al Sallab, Senthil Yogamani, and Patrick Pérez. Deep reinforcement learning for autonomous driving: A survey. <i>IEEE</i> <i>Transactions on Intelligent Transportation Systems</i>, 23(6):4909–4926, 2021. Ilya Kostrikov, Denis Yarats, and Rob Fergus. Image augmentation is all you need: Regularizing deep reinforcement learning from pixels. <i>arXiv preprint arXiv:2004.13649</i>, 2020. Ilya Kostrikov, Ashvin Nair, and Sergey Levine. Offline reinforcement learning with implicit Q-learning. <i>arXiv preprint arXiv:2110.06169</i>, 2021. Aviral Kumar, Justin Fu, Matthew Soh, George Tucker, and Sergey Levine. Stabilizing off-policy q-learning via bootstrapping error reduction. <i>Advances in Neural Information Processing Systems</i>, 32, 2019. Michael Laskin, Aravind Srinivas, and Pieter Abbeel. CURL: Contrastive unsupervised representa- tions for reinforcement learning. In <i>International Conference on Machine Learning</i>, pp. 5639–5650. PMLR, 2020. Charline Le Lan, Marc G Bellemare, and Pablo Samuel Castro. Metrics and continuity in reinforce- ment learning. In <i>Proceedings of the AAAI Conference on Artificial Intelligence</i>, volume 35, pp. 8261–8269, 2021. Jongmin Lee, Wonseok Jeon, Byungjun Lee, Joelle Pineau, and Kee-Eung Kim. OptiDICE: Offline policy optimization via stationary distribution correction estimation. In <i>International Conference on Machine Learning</i>, pp. 6120–6130. PMLR, 2021. Jongmin Lee, Cosmin Paduraru, Daniel J Mankowitz, Nicolas Heess, Doina Precup, Kee-Eung Kim, and Arthur Guez. COptiDICE: Offline constrained reinforcement learning via stationary distribution correction estimation. <i>arXiv preprint arXiv:2204.08957</i>, 2022. Vint Lee, Pieter Abbeel, and Youngwoon Lee. Dreamsmooth: Improving model-based reinforcement learning via reward smoothing. In <i>The Twelfth International Conference on Learning Representa-</i> <i>tions 2024</i> 	620	
 Transactions on Intelligent Transportation Systems, 23(6):4909–4926, 2021. Ilya Kostrikov, Denis Yarats, and Rob Fergus. Image augmentation is all you need: Regularizing deep reinforcement learning from pixels. <i>arXiv preprint arXiv:2004.13649</i>, 2020. Ilya Kostrikov, Ashvin Nair, and Sergey Levine. Offline reinforcement learning with implicit Q-learning. <i>arXiv preprint arXiv:2110.06169</i>, 2021. Aviral Kumar, Justin Fu, Matthew Soh, George Tucker, and Sergey Levine. Stabilizing off-policy q-learning via bootstrapping error reduction. <i>Advances in Neural Information Processing Systems</i>, 32, 2019. Michael Laskin, Aravind Srinivas, and Pieter Abbeel. CURL: Contrastive unsupervised representations for reinforcement learning. In <i>International Conference on Machine Learning</i>, pp. 5639–5650. PMLR, 2020. Charline Le Lan, Marc G Bellemare, and Pablo Samuel Castro. Metrics and continuity in reinforcement learning. In <i>Proceedings of the AAAI Conference on Artificial Intelligence</i>, volume 35, pp. 8261–8269, 2021. Jongmin Lee, Wonseok Jeon, Byungjun Lee, Joelle Pineau, and Kee-Eung Kim. OptiDICE: Offline policy optimization via stationary distribution correction estimation. In <i>International Conference on Machine Learning</i>, pp. 6120–6130. PMLR, 2021. Jongmin Lee, Cosmin Paduraru, Daniel J Mankowitz, Nicolas Heess, Doina Precup, Kee-Eung Kim, and Arthur Guez. COptiDICE: Offline constrained reinforcement learning via stationary distribution correction estimation arXiv <i>preprint arXiv:2204.08957</i>, 2022. Vint Lee, Pieter Abbeel, and Youngwoon Lee. Dreamsmooth: Improving model-based reinforcement learning via reward smoothing. In <i>The Twelfth International Conference on Learning Representations</i>, 2024 	621 622	B Ravi Kiran, Ibrahim Sobh, Victor Talpaert, Patrick Mannion, Ahmad A Al Sallab, Senthil Yogamani, and Patrick Pérez. Deep reinforcement learning for autonomous driving: A survey. <i>IEEE</i>
 Ilya Kostrikov, Denis Yarats, and Rob Fergus. Image augmentation is all you need: Regularizing deep reinforcement learning from pixels. <i>arXiv preprint arXiv:2004.13649</i>, 2020. Ilya Kostrikov, Ashvin Nair, and Sergey Levine. Offline reinforcement learning with implicit Q-learning. <i>arXiv preprint arXiv:2110.06169</i>, 2021. Aviral Kumar, Justin Fu, Matthew Soh, George Tucker, and Sergey Levine. Stabilizing off-policy q-learning via bootstrapping error reduction. <i>Advances in Neural Information Processing Systems</i>, 32, 2019. Michael Laskin, Aravind Srinivas, and Pieter Abbeel. CURL: Contrastive unsupervised representations for reinforcement learning. In <i>International Conference on Machine Learning</i>, pp. 5639–5650. PMLR, 2020. Charline Le Lan, Marc G Bellemare, and Pablo Samuel Castro. Metrics and continuity in reinforcement learning. In <i>Proceedings of the AAAI Conference on Artificial Intelligence</i>, volume 35, pp. 8261–8269, 2021. Jongmin Lee, Wonseok Jeon, Byungjun Lee, Joelle Pineau, and Kee-Eung Kim. OptiDICE: Offline policy optimization via stationary distribution correction estimation. In <i>International Conference on Machine Learning</i>. Nachura, pp. 6120–6130. PMLR, 2021. Jongmin Lee, Cosmin Paduraru, Daniel J Mankowitz, Nicolas Heess, Doina Precup, Kee-Eung Kim, and Arthur Guez. COptiDICE: Offline constrained reinforcement learning via stationary distribution correction estimation arXiv:2204.08957, 2022. Vint Lee, Pieter Abbeel, and Youngwoon Lee. Dreamsmooth: Improving model-based reinforcement learning reward smoothing. In <i>The Welfth International Conference on Learning Representations</i>, 2024. 	623	Transactions on Intelligent Transportation Systems, 23(6):4909–4926, 2021.
 Ilya Kostrikov, Ashvin Nair, and Sergey Levine. Offline reinforcement learning with implicit Q-learning. <i>arXiv preprint arXiv:2110.06169</i>, 2021. Aviral Kumar, Justin Fu, Matthew Soh, George Tucker, and Sergey Levine. Stabilizing off-policy q-learning via bootstrapping error reduction. <i>Advances in Neural Information Processing Systems</i>, 32, 2019. Michael Laskin, Aravind Srinivas, and Pieter Abbeel. CURL: Contrastive unsupervised representa- tions for reinforcement learning. In <i>International Conference on Machine Learning</i>, pp. 5639–5650. PMLR, 2020. Charline Le Lan, Marc G Bellemare, and Pablo Samuel Castro. Metrics and continuity in reinforce- ment learning. In <i>Proceedings of the AAAI Conference on Artificial Intelligence</i>, volume 35, pp. 8261–8269, 2021. Jongmin Lee, Wonseok Jeon, Byungjun Lee, Joelle Pineau, and Kee-Eung Kim. OptiDICE: Offline policy optimization via stationary distribution correction estimation. In <i>International Conference on Machine Learning</i>, pp. 6120–6130. PMLR, 2021. Jongmin Lee, Cosmin Paduraru, Daniel J Mankowitz, Nicolas Heess, Doina Precup, Kee-Eung Kim, and Arthur Guez. COptiDICE: Offline constrained reinforcement learning via stationary distribution correction estimation. <i>arXiv preprint arXiv:2204.08957</i>, 2022. Vint Lee, Pieter Abbeel, and Youngwoon Lee. Dreamsmooth: Improving model-based reinforcement learning via reward smoothing. In <i>The Twelfth International Conference on Learning Representa- tions</i>, 2024 	624 625	Ilya Kostrikov, Denis Yarats, and Rob Fergus. Image augmentation is all you need: Regularizing deep reinforcement learning from pixels. <i>arXiv preprint arXiv:2004.13649</i> , 2020.
 Aviral Kumar, Justin Fu, Matthew Soh, George Tucker, and Sergey Levine. Stabilizing off-policy q-learning via bootstrapping error reduction. <i>Advances in Neural Information Processing Systems</i>, 32, 2019. Michael Laskin, Aravind Srinivas, and Pieter Abbeel. CURL: Contrastive unsupervised representa- tions for reinforcement learning. In <i>International Conference on Machine Learning</i>, pp. 5639–5650. PMLR, 2020. Charline Le Lan, Marc G Bellemare, and Pablo Samuel Castro. Metrics and continuity in reinforce- ment learning. In <i>Proceedings of the AAAI Conference on Artificial Intelligence</i>, volume 35, pp. 8261–8269, 2021. Jongmin Lee, Wonseok Jeon, Byungjun Lee, Joelle Pineau, and Kee-Eung Kim. OptiDICE: Offline policy optimization via stationary distribution correction estimation. In <i>International Conference on Machine Learning</i>, pp. 6120–6130. PMLR, 2021. Jongmin Lee, Cosmin Paduraru, Daniel J Mankowitz, Nicolas Heess, Doina Precup, Kee-Eung Kim, and Arthur Guez. COptiDICE: Offline constrained reinforcement learning via stationary distribution correction estimation. <i>arXiv preprint arXiv:2204.08957</i>, 2022. Vint Lee, Pieter Abbeel, and Youngwoon Lee. Dreamsmooth: Improving model-based reinforcement learning via reward smoothing. In <i>The Twelfth International Conference on Learning Representa- tions</i>, 2024. 	626 627	Ilya Kostrikov, Ashvin Nair, and Sergey Levine. Offline reinforcement learning with implicit Q-learning. <i>arXiv preprint arXiv:2110.06169</i> , 2021.
 Aviral Kullia, Justin Fu, Matthew Son, George Tucker, and Sergey Levine. Stabilizing on-policy q-learning via bootstrapping error reduction. Advances in Neural Information Processing Systems, 32, 2019. Michael Laskin, Aravind Srinivas, and Pieter Abbeel. CURL: Contrastive unsupervised representa- tions for reinforcement learning. In International Conference on Machine Learning, pp. 5639–5650. PMLR, 2020. Charline Le Lan, Marc G Bellemare, and Pablo Samuel Castro. Metrics and continuity in reinforce- ment learning. In Proceedings of the AAAI Conference on Artificial Intelligence, volume 35, pp. 8261–8269, 2021. Jongmin Lee, Wonseok Jeon, Byungjun Lee, Joelle Pineau, and Kee-Eung Kim. OptiDICE: Offline policy optimization via stationary distribution correction estimation. In International Conference on Machine Learning, pp. 6120–6130. PMLR, 2021. Jongmin Lee, Cosmin Paduraru, Daniel J Mankowitz, Nicolas Heess, Doina Precup, Kee-Eung Kim, and Arthur Guez. COptiDICE: Offline constrained reinforcement learning via stationary distribution correction estimation. arXiv preprint arXiv:2204.08957, 2022. Vint Lee, Pieter Abbeel, and Youngwoon Lee. Dreamsmooth: Improving model-based reinforcement learning via reward smoothing. In The Twelfth International Conference on Learning Representa- tions 2024 	620	Avies Kymen Lystin Fy Matthew Sch. Coorse Typker and Sergery Levine Stabilizing off policy
 Michael Laskin, Aravind Srinivas, and Pieter Abbeel. CURL: Contrastive unsupervised representations for reinforcement learning. In <i>International Conference on Machine Learning</i>, pp. 5639–5650. PMLR, 2020. Charline Le Lan, Marc G Bellemare, and Pablo Samuel Castro. Metrics and continuity in reinforcement learning. In <i>Proceedings of the AAAI Conference on Artificial Intelligence</i>, volume 35, pp. 8261–8269, 2021. Jongmin Lee, Wonseok Jeon, Byungjun Lee, Joelle Pineau, and Kee-Eung Kim. OptiDICE: Offline policy optimization via stationary distribution correction estimation. In <i>International Conference on Machine Learning</i>, pp. 6120–6130. PMLR, 2021. Jongmin Lee, Cosmin Paduraru, Daniel J Mankowitz, Nicolas Heess, Doina Precup, Kee-Eung Kim, and Arthur Guez. COptiDICE: Offline constrained reinforcement learning via stationary distribution correction estimation. <i>arXiv preprint arXiv:2204.08957</i>, 2022. Vint Lee, Pieter Abbeel, and Youngwoon Lee. Dreamsmooth: Improving model-based reinforcement learning reinforcement learning Representations 2024 	629 630	q-learning via bootstrapping error reduction. <i>Advances in Neural Information Processing Systems</i> , 32, 2019
 Michael Laskin, Aravind Srinivas, and Pieter Abbeel. CURL: Contrastive unsupervised representations for reinforcement learning. In <i>International Conference on Machine Learning</i>, pp. 5639–5650. PMLR, 2020. Charline Le Lan, Marc G Bellemare, and Pablo Samuel Castro. Metrics and continuity in reinforcement learning. In <i>Proceedings of the AAAI Conference on Artificial Intelligence</i>, volume 35, pp. 8261–8269, 2021. Jongmin Lee, Wonseok Jeon, Byungjun Lee, Joelle Pineau, and Kee-Eung Kim. OptiDICE: Offline policy optimization via stationary distribution correction estimation. In <i>International Conference on Machine Learning</i>, pp. 6120–6130. PMLR, 2021. Jongmin Lee, Cosmin Paduraru, Daniel J Mankowitz, Nicolas Heess, Doina Precup, Kee-Eung Kim, and Arthur Guez. COptiDICE: Offline constrained reinforcement learning via stationary distribution correction estimation. arXiv preprint arXiv:2204.08957, 2022. Vint Lee, Pieter Abbeel, and Youngwoon Lee. Dreamsmooth: Improving model-based reinforcement learning via reward smoothing. In <i>The Twelfth International Conference on Learning Representations</i> 2024. 	031	52, 2019.
 tions for reinforcement learning. In <i>International Conference on Machine Learning</i>, pp. 5639–5650. PMLR, 2020. Charline Le Lan, Marc G Bellemare, and Pablo Samuel Castro. Metrics and continuity in reinforcement learning. In <i>Proceedings of the AAAI Conference on Artificial Intelligence</i>, volume 35, pp. 8261–8269, 2021. Jongmin Lee, Wonseok Jeon, Byungjun Lee, Joelle Pineau, and Kee-Eung Kim. OptiDICE: Offline policy optimization via stationary distribution correction estimation. In <i>International Conference on Machine Learning</i>, pp. 6120–6130. PMLR, 2021. Jongmin Lee, Cosmin Paduraru, Daniel J Mankowitz, Nicolas Heess, Doina Precup, Kee-Eung Kim, and Arthur Guez. COptiDICE: Offline constrained reinforcement learning via stationary distribution correction estimation. <i>arXiv preprint arXiv:2204.08957</i>, 2022. Vint Lee, Pieter Abbeel, and Youngwoon Lee. Dreamsmooth: Improving model-based reinforcement learning via reward smoothing. In <i>The Twelfth International Conference on Learning Representations</i>, 2024. 	632	Michael Laskin, Aravind Srinivas, and Pieter Abbeel. CURL: Contrastive unsupervised representa-
 PMLR, 2020. Charline Le Lan, Marc G Bellemare, and Pablo Samuel Castro. Metrics and continuity in reinforce- ment learning. In <i>Proceedings of the AAAI Conference on Artificial Intelligence</i>, volume 35, pp. 8261–8269, 2021. Jongmin Lee, Wonseok Jeon, Byungjun Lee, Joelle Pineau, and Kee-Eung Kim. OptiDICE: Offline policy optimization via stationary distribution correction estimation. In <i>International Conference on Machine Learning</i>, pp. 6120–6130. PMLR, 2021. Jongmin Lee, Cosmin Paduraru, Daniel J Mankowitz, Nicolas Heess, Doina Precup, Kee-Eung Kim, and Arthur Guez. COptiDICE: Offline constrained reinforcement learning via stationary distribution correction estimation. <i>arXiv preprint arXiv:2204.08957</i>, 2022. Vint Lee, Pieter Abbeel, and Youngwoon Lee. Dreamsmooth: Improving model-based reinforcement learning via reward smoothing. In <i>The Twelfth International Conference on Learning Representations</i>, 2024. 	634	tions for reinforcement learning. In International Conference on Machine Learning, pp. 5639–5650.
 Charline Le Lan, Marc G Bellemare, and Pablo Samuel Castro. Metrics and continuity in reinforce- ment learning. In <i>Proceedings of the AAAI Conference on Artificial Intelligence</i>, volume 35, pp. 8261–8269, 2021. Jongmin Lee, Wonseok Jeon, Byungjun Lee, Joelle Pineau, and Kee-Eung Kim. OptiDICE: Offline policy optimization via stationary distribution correction estimation. In <i>International Conference on Machine Learning</i>, pp. 6120–6130. PMLR, 2021. Jongmin Lee, Cosmin Paduraru, Daniel J Mankowitz, Nicolas Heess, Doina Precup, Kee-Eung Kim, and Arthur Guez. COptiDICE: Offline constrained reinforcement learning via stationary distribution correction estimation. <i>arXiv preprint arXiv:2204.08957</i>, 2022. Vint Lee, Pieter Abbeel, and Youngwoon Lee. Dreamsmooth: Improving model-based reinforcement learning via reward smoothing. In <i>The Twelfth International Conference on Learning Representations</i>, 2024. 	625	PMLR, 2020.
 ment learning. In <i>Proceedings of the AAAI Conference on Artificial Intelligence</i>, volume 35, pp. 8261–8269, 2021. Jongmin Lee, Wonseok Jeon, Byungjun Lee, Joelle Pineau, and Kee-Eung Kim. OptiDICE: Offline policy optimization via stationary distribution correction estimation. In <i>International Conference on Machine Learning</i>, pp. 6120–6130. PMLR, 2021. Jongmin Lee, Cosmin Paduraru, Daniel J Mankowitz, Nicolas Heess, Doina Precup, Kee-Eung Kim, and Arthur Guez. COptiDICE: Offline constrained reinforcement learning via stationary distribution correction estimation. <i>arXiv preprint arXiv:2204.08957</i>, 2022. Vint Lee, Pieter Abbeel, and Youngwoon Lee. Dreamsmooth: Improving model-based reinforcement learning via reward smoothing. In <i>The Twelfth International Conference on Learning Representations</i>, 2024. 	636	Charline Le Lan, Marc G Bellemare, and Pablo Samuel Castro. Metrics and continuity in reinforce-
 8261–8269, 2021. Jongmin Lee, Wonseok Jeon, Byungjun Lee, Joelle Pineau, and Kee-Eung Kim. OptiDICE: Offline policy optimization via stationary distribution correction estimation. In <i>International Conference</i> <i>on Machine Learning</i>, pp. 6120–6130. PMLR, 2021. Jongmin Lee, Cosmin Paduraru, Daniel J Mankowitz, Nicolas Heess, Doina Precup, Kee-Eung Kim, and Arthur Guez. COptiDICE: Offline constrained reinforcement learning via stationary distribution correction estimation. <i>arXiv preprint arXiv:2204.08957</i>, 2022. Vint Lee, Pieter Abbeel, and Youngwoon Lee. Dreamsmooth: Improving model-based reinforcement learning via reward smoothing. In <i>The Twelfth International Conference on Learning Representa-</i> 	637	ment learning. In Proceedings of the AAAI Conference on Artificial Intelligence, volume 35, pp.
 Jongmin Lee, Wonseok Jeon, Byungjun Lee, Joelle Pineau, and Kee-Eung Kim. OptiDICE: Offline policy optimization via stationary distribution correction estimation. In <i>International Conference</i> <i>on Machine Learning</i>, pp. 6120–6130. PMLR, 2021. Jongmin Lee, Cosmin Paduraru, Daniel J Mankowitz, Nicolas Heess, Doina Precup, Kee-Eung Kim, and Arthur Guez. COptiDICE: Offline constrained reinforcement learning via stationary distribution correction estimation. <i>arXiv preprint arXiv:2204.08957</i>, 2022. Vint Lee, Pieter Abbeel, and Youngwoon Lee. Dreamsmooth: Improving model-based reinforcement learning via reward smoothing. In <i>The Twelfth International Conference on Learning Representa-</i> 	638	8261–8269, 2021.
 Jongmin Lee, wonseok Jeon, Byungjun Lee, Joene Pineau, and Ree-Eding Kini. OptiDICE. Online policy optimization via stationary distribution correction estimation. In <i>International Conference</i> <i>on Machine Learning</i>, pp. 6120–6130. PMLR, 2021. Jongmin Lee, Cosmin Paduraru, Daniel J Mankowitz, Nicolas Heess, Doina Precup, Kee-Eung Kim, and Arthur Guez. COptiDICE: Offline constrained reinforcement learning via stationary distribution correction estimation. <i>arXiv preprint arXiv:2204.08957</i>, 2022. Vint Lee, Pieter Abbeel, and Youngwoon Lee. Dreamsmooth: Improving model-based reinforcement learning via reward smoothing. In <i>The Twelfth International Conference on Learning Representa-</i> 	639	Jongmin Los Wonsook Joon Puungiun Los Joelle Pineeu, and Koo Fung Kim. OntiDICE: Offline
 641 on Machine Learning, pp. 6120–6130. PMLR, 2021. 642 643 Jongmin Lee, Cosmin Paduraru, Daniel J Mankowitz, Nicolas Heess, Doina Precup, Kee-Eung 644 Kim, and Arthur Guez. COptiDICE: Offline constrained reinforcement learning via stationary 645 distribution correction estimation. arXiv preprint arXiv:2204.08957, 2022. 646 Vint Lee, Pieter Abbeel, and Youngwoon Lee. Dreamsmooth: Improving model-based reinforcement 647 learning via reward smoothing. In <i>The Twelfth International Conference on Learning Representa-</i> 647 <i>tions</i> 2024 	640	policy optimization via stationary distribution correction estimation. In International Conference
 Jongmin Lee, Cosmin Paduraru, Daniel J Mankowitz, Nicolas Heess, Doina Precup, Kee-Eung Kim, and Arthur Guez. COptiDICE: Offline constrained reinforcement learning via stationary distribution correction estimation. <i>arXiv preprint arXiv:2204.08957</i>, 2022. Vint Lee, Pieter Abbeel, and Youngwoon Lee. Dreamsmooth: Improving model-based reinforcement learning via reward smoothing. In <i>The Twelfth International Conference on Learning Representa-</i> <i>tions</i>, 2024. 	641	on Machine Learning, pp. 6120–6130. PMLR, 2021.
 Kim, and Arthur Guez. COptiDICE: Offline constrained reinforcement learning via stationary distribution correction estimation. <i>arXiv preprint arXiv:2204.08957</i>, 2022. Vint Lee, Pieter Abbeel, and Youngwoon Lee. Dreamsmooth: Improving model-based reinforcement learning via reward smoothing. In <i>The Twelfth International Conference on Learning Representations</i> 2024 	642	Iongmin Lee, Cosmin Paduraru, Daniel I Mankowitz, Nicolas Heess, Doina Precup, Kee-Fung
 distribution correction estimation. arXiv preprint arXiv:2204.08957, 2022. Vint Lee, Pieter Abbeel, and Youngwoon Lee. Dreamsmooth: Improving model-based reinforcement learning via reward smoothing. In <i>The Twelfth International Conference on Learning Representa-</i> <i>tions</i>, 2024. 	643	Kim, and Arthur Guez. COptiDICE: Offline constrained reinforcement learning via stationary
 645 646 Vint Lee, Pieter Abbeel, and Youngwoon Lee. Dreamsmooth: Improving model-based reinforcement 647 learning via reward smoothing. In <i>The Twelfth International Conference on Learning Representa-</i> <i>tions</i> 2024 	644	distribution correction estimation. arXiv preprint arXiv:2204.08957, 2022.
 646 Vint Lee, Pieter Abbeel, and Youngwoon Lee. Dreamsmooth: Improving model-based reinforcement 647 learning via reward smoothing. In <i>The Twelfth International Conference on Learning Representa-</i> tions 2024 	645	
	646 647	Vint Lee, Pieter Abbeel, and Youngwoon Lee. Dreamsmooth: Improving model-based reinforcement learning via reward smoothing. In <i>The Twelfth International Conference on Learning Representa-</i> <i>tions</i> 2024

652

657

- 648 Lihong Li, Thomas J Walsh, and Michael L Littman. Towards a unified theory of state abstraction for 649 MDPs. AI&M, 1(2):3, 2006. 650
 - Shengbo Eben Li. Reinforcement learning for sequential decision and optimal control. Springer, 2023.
- 653 Qian Lin, Bo Tang, Zifan Wu, Chao Yu, Shangqin Mao, Qianlong Xie, Xingxing Wang, and Dong 654 Wang. Safe offline reinforcement learning with real-time budget constraints. In International 655 Conference on Machine Learning, pp. 21127–21152. PMLR, 2023.
- 656 Guoqing Liu, Chuheng Zhang, Li Zhao, Tao Qin, Jinhua Zhu, Jian Li, Nenghai Yu, and Tie-Yan Liu. Return-based contrastive representation learning for reinforcement learning. arXiv preprint 658 arXiv:2102.10960, 2021. 659
- Zuxin Liu, Zhepeng Cen, Vladislav Isenbaev, Wei Liu, Steven Wu, Bo Li, and Ding Zhao. Constrained 660 variational policy optimization for safe reinforcement learning. In International Conference on 661 Machine Learning, pp. 13644–13668. PMLR, 2022. 662
- 663 Zuxin Liu, Zijian Guo, Haohong Lin, Yihang Yao, Jiacheng Zhu, Zhepeng Cen, Hanjiang Hu, Wenhao 664 Yu, Tingnan Zhang, Jie Tan, et al. Datasets and benchmarks for offline safe reinforcement learning. 665 *arXiv preprint arXiv:2306.09303*, 2023a.
- 666 Zuxin Liu, Zijian Guo, Yihang Yao, Zhepeng Cen, Wenhao Yu, Tingnan Zhang, and Ding Zhao. Con-667 strained decision transformer for offline safe reinforcement learning. In International Conference 668 on Machine Learning, pp. 21611–21630. PMLR, 2023b. 669
- Cheng Lu, Yuhao Zhou, Fan Bao, Jianfei Chen, Chongxuan Li, and Jun Zhu. DPM-solver: A fast 670 ode solver for diffusion probabilistic model sampling in around 10 steps. Advances in Neural 671 Information Processing Systems, 35:5775–5787, 2022a. 672
- 673 Cheng Lu, Yuhao Zhou, Fan Bao, Jianfei Chen, Chongxuan Li, and Jun Zhu. DPM-solver++: Fast 674 solver for guided sampling of diffusion probabilistic models. arXiv preprint arXiv:2211.01095, 675 2022b.
- 676 Cheng Lu, Huayu Chen, Jianfei Chen, Hang Su, Chongxuan Li, and Jun Zhu. Contrastive energy 677 prediction for exact energy-guided diffusion sampling in offline reinforcement learning. In 678 International Conference on Machine Learning, pp. 22825–22855. PMLR, 2023. 679
- Yuping Luo and Tengyu Ma. Learning barrier certificates: Towards safe reinforcement learning 680 with zero training-time violations. Advances in Neural Information Processing Systems, 34: 681 25621-25632, 2021. 682
- 683 Haitong Ma, Jianyu Chen, Shengbo Eben, Ziyu Lin, Yang Guan, Yangang Ren, and Sifa Zheng. 684 Model-based constrained reinforcement learning using generalized control barrier function. In 2021 685 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS), pp. 4552–4559. IEEE, 2021a. 686
- 687 Haitong Ma, Yang Guan, Shegnbo Eben Li, Xiangteng Zhang, Sifa Zheng, and Jianyu Chen. Feasible 688 actor-critic: Constrained reinforcement learning for ensuring statewise safety. arXiv preprint 689 arXiv:2105.10682, 2021b. 690
- Bhairav Mehta, Tristan Deleu, Sharath Chandra Raparthy, Chris J Pal, and Liam Paull. Curriculum in 691 gradient-based meta-reinforcement learning. arXiv preprint arXiv:2002.07956, 2020. 692
- 693 Volodymyr Mnih, Koray Kavukcuoglu, David Silver, Andrei A Rusu, Joel Veness, Marc G Bellemare, 694 Alex Graves, Martin Riedmiller, Andreas K Fidjeland, Georg Ostrovski, et al. Human-level control 695 through deep reinforcement learning. *Nature*, 518(7540):529–533, 2015.
- 696 Aaron van den Oord, Yazhe Li, and Oriol Vinyals. Representation learning with contrastive predictive 697 coding. arXiv preprint arXiv:1807.03748, 2018. 698
- Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor 699 Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, et al. Pytorch: An imperative style, 700 high-performance deep learning library. Advances in Neural Information Processing Systems, 32, 701 2019.

- 702 Zhenghao Peng, Quanyi Li, Chunxiao Liu, and Bolei Zhou. Safe driving via expert guided policy 703 optimization. In Conference on Robot Learning, pp. 1554–1563. PMLR, 2022. 704 705 Alex Ray, Joshua Achiam, and Dario Amodei. Benchmarking safe exploration in deep reinforcement learning. arXiv preprint arXiv:1910.01708, 7(1):2, 2019. 706 707 Julian Schrittwieser, Ioannis Antonoglou, Thomas Hubert, Karen Simonyan, Laurent Sifre, Simon 708 Schmitt, Arthur Guez, Edward Lockhart, Demis Hassabis, Thore Graepel, et al. Mastering Atari, 709 Go, chess and shogi by planning with a learned model. Nature, 588(7839):604-609, 2020. 710 David Silver, Julian Schrittwieser, Karen Simonyan, Ioannis Antonoglou, Aja Huang, Arthur Guez, 711 Thomas Hubert, Lucas Baker, Matthew Lai, Adrian Bolton, et al. Mastering the game of Go 712 without human knowledge. Nature, 550(7676):354-359, 2017. 713 714 Jascha Sohl-Dickstein, Eric Weiss, Niru Maheswaranathan, and Surya Ganguli. Deep unsupervised 715 learning using nonequilibrium thermodynamics. In International Conference on Machine Learning, 716 pp. 2256–2265. PMLR, 2015. 717 Yang Song, Jascha Sohl-Dickstein, Diederik P Kingma, Abhishek Kumar, Stefano Ermon, and Ben 718 Poole. Score-based generative modeling through stochastic differential equations. arXiv preprint 719 arXiv:2011.13456, 2020. 720 Aivar Sootla, Alexander I Cowen-Rivers, Taher Jafferjee, Ziyan Wang, David H Mguni, Jun Wang, and 721 Haitham Ammar. Sauté RL: Almost surely safe reinforcement learning using state augmentation. 722 In International Conference on Machine Learning, pp. 20423–20443. PMLR, 2022. 723 724 Adam Stooke, Joshua Achiam, and Pieter Abbeel. Responsive safety in reinforcement learning by pid 725 lagrangian methods. In International Conference on Machine Learning, pp. 9133–9143. PMLR, 726 2020. 727 Jayakumar Subramanian, Amit Sinha, Raihan Seraj, and Aditya Mahajan. Approximate information 728 state for approximate planning and reinforcement learning in partially observed systems. Journal 729 of Machine Learning Research, 23(12):1-83, 2022. 730 731 Chen Tessler, Daniel J Mankowitz, and Shie Mannor. Reward constrained policy optimization. arXiv preprint arXiv:1805.11074, 2018. 732 733 Brijen Thananjeyan, Ashwin Balakrishna, Suraj Nair, Michael Luo, Krishnan Srinivasan, Minho 734 Hwang, Joseph E Gonzalez, Julian Ibarz, Chelsea Finn, and Ken Goldberg. Recovery RL: Safe 735 reinforcement learning with learned recovery zones. IEEE Robotics and Automation Letters, 6(3): 736 4915-4922, 2021. 737 Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz 738 Kaiser, and Illia Polosukhin. Attention is all you need. Advances in Neural Information Processing 739 Systems, 30, 2017. 740 741 Akifumi Wachi, Yanan Sui, Yisong Yue, and Masahiro Ono. Safe exploration and optimization of 742 constrained mdps using gaussian processes. In Proceedings of the AAAI Conference on Artificial 743 Intelligence, volume 32, 2018. 744 Zhendong Wang, Jonathan J Hunt, and Mingyuan Zhou. Diffusion policies as an expressive policy 745 class for offline reinforcement learning. arXiv preprint arXiv:2208.06193, 2022. 746 747 Manuel Watter, Jost Springenberg, Joschka Boedecker, and Martin Riedmiller. Embed to control: A 748 locally linear latent dynamics model for control from raw images. Advances in Neural Information Processing Systems, 28, 2015. 749 750 Haoran Xu, Xianyuan Zhan, and Xiangyu Zhu. Constraints penalized Q-learning for safe offline rein-751 forcement learning. In Proceedings of the AAAI Conference on Artificial Intelligence, volume 36, 752 pp. 8753-8760, 2022a. 753 Haoran Xu, Li Jiang, Jianxiong Li, Zhuoran Yang, Zhaoran Wang, Victor Wai Kin Chan, and Xianyuan 754
- 755 Zhao, Offline RL with no ood actions: In-sample learning via implicit value regularization. *arXiv* preprint arXiv:2303.15810, 2023.

756 757 758	Mengdi Xu, Zuxin Liu, Peide Huang, Wenhao Ding, Zhepeng Cen, Bo Li, and Ding Zhao. Trustworthy reinforcement learning against intrinsic vulnerabilities: Robustness, safety, and generalizability. <i>arXiv preprint arXiv:2209.08025</i> , 2022b.
759 760 761	Zhenghai Xue, Zhenghao Peng, Quanyi Li, Zhihan Liu, and Bolei Zhou. Guarded policy optimization with imperfect online demonstrations. <i>arXiv preprint arXiv:2303.01728</i> , 2023.
762 763 764	Denis Yarats, Amy Zhang, Ilya Kostrikov, Brandon Amos, Joelle Pineau, and Rob Fergus. Improving sample efficiency in model-free reinforcement learning from images. In <i>Proceedings of the AAAI Conference on Artificial Intelligence</i> , volume 35, pp. 10674–10681, 2021.
765 766 767	Dongjie Yu, Haitong Ma, Shengbo Li, and Jianyu Chen. Reachability constrained reinforcement learning. In <i>International Conference on Machine Learning</i> , pp. 25636–25655. PMLR, 2022.
768 769	Weiye Zhao, Tairan He, and Changliu Liu. Model-free safe control for zero-violation reinforcement learning. In 5th Annual Conference on Robot Learning, 2021.
770 771 772 773	Yinan Zheng, Jianxiong Li, Dongjie Yu, Yujie Yang, Shengbo Eben Li, Xianyuan Zhan, and Jingjing Liu. Safe offline reinforcement learning with feasibility-guided diffusion model. <i>arXiv preprint arXiv:2401.10700</i> , 2024.
774 775 776	Jinhua Zhu, Yingce Xia, Lijun Wu, Jiajun Deng, Wengang Zhou, Tao Qin, Tie-Yan Liu, and Houqiang Li. Masked contrastive representation learning for reinforcement learning. <i>IEEE Transactions on Pattern Analysis and Machine Intelligence</i> , 45(3):3421–3433, 2022.
777	
778	
779	
780	
781	
782	
783	
784	
785	
786	
700	
780	
705	
791	
792	
793	
794	
795	
796	
797	
798	
799	
800	
801	
802	
803	
804	
805	
806	
007	
809	
003	

810 THEORETICAL INTERPRETATIONS А 811

812 The first comparison between the bisimulation representation and Q^* -irrelevance representation for 813 finite horizon MDPs was given by Givan et al. (2003). The systematic state abstraction theory for 814 MDPs was summarized in Li et al. (2006). The expansion of the theory in Partially Observable MDPs 815 (POMDPs) are introduced recently (Subramanian et al., 2022). It is worth noting that their formulation 816 does not incorporate the safety Bellman operator, and a comprehensive proof for infinite-horizon MDPs is not provided. We now provide a complete proof for Theorem 3.1 as follows. 817

818 **Definition A.1.** For any state representation Θ_1, Θ_2 , we say Θ_1 is finer than Θ_2 , denoted as $\Theta_1 \succeq \Theta_2$, if and only if for any states $s_1, s_2 \in S$, $\Theta_1(s_1) = \Theta_1(s_2)$ implies $\Theta_2(s_1) = \Theta_2(s_2)$. 819

To clarify, let z_1 and z_2 represent the representations $\Theta_1(s)$ and $\Theta_2(s)$ for any $s \in S$, respectively. It 821 is always possible to find a function $f: \mathcal{Z}_1 \to \mathcal{Z}_2$ that is surjective. The equality holds $(\Theta_1 = \Theta_2)$ if 822 and only if the surjective function is also injective (i.e., bijective). 823

Theorem 3.1. For any MDP, the optimal Q-functions induced by either the general Bellman operator 824 or the safety Bellman operator satisfy $\Theta_{\text{bisim}} \succeq \Theta_{Q^*}$. The optimal policies derived from both 825 bisimulation representation and Q^* -irrelevance representation are also optimal in the ground MDP, 826 *i.e.*, $\pi^*(\Theta_{\text{bisim}}(s)) = \pi^*(\Theta_{Q^*}(s)) = \pi^*(s)$ for any state $s \in S$. 827

828 *Proof.* We start by considering a finite-horizon MDP with a maximum timestep T. For any timestep 829 $t \in \{1, 2, ..., T\}$, we denote $Q_{r,t(T)}^*$ as the optimal-Q function at timestep t. Then, for $\forall s \in S$ and 830 $\forall a \in \mathcal{A}$, we have: 831

820

$$Q_{r,t(T)}^{*}(s,a) = r(s,a) + \gamma \sum_{s' \in \mathcal{S}} P(s' \mid s_1, a) [\max_{a' \in \mathcal{A}} Q_{r,t+1(T)}^{*}(s',a')].$$
(16)

833 For timestep T + 1, we define $Q_{r,T+1(T)}^*(s,a) = 0$ for $\forall s \in S$ and $\forall a \in A$, which implies 834 that $Q^*_{r,T(T)}(s,a) = r(s,a)$. Now, for any $s_1, s_2 \in S$ that are bisimilar (i.e., $\Theta_{bisim}(s_1) =$ 835 $\Theta_{bisim}(s_2)$), we have $Q^*_{r,T(T)}(s_1,a) = Q^*_{r,T(T)}(s_2,a)$. In other words, for any $z' \in \mathbb{Z}_{bisim}$, 836 $\max_{a' \in \mathcal{A}} Q_{r,T(T)}^*(s',a')$ is the same for all $s' \in \Theta_{bisim}^{-1}(z')$. 837

838 Considering any $s_1, s_2 \in S$ that are bisimilar, and for any action $a \in A$, we perform backward 839 induction on timestep t from T-1 to 1 following the proof sketch of Theorem 5 in Givan et al. 840 (2003): 841

842 843

844 845

846

847

848 849

$$= r(s_{1}, a) + \gamma \sum_{s' \in \mathcal{S}} P(s' \mid s_{1}, a) [\max_{a' \in \mathcal{A}} Q_{r,t+1(T)}^{*}(s', a')]$$

$$\stackrel{a}{=} r(s_{1}, a) + \gamma \sum_{s' \in \{\bigcup_{z' \in \mathcal{Z}} \Theta^{-1}(z')\}} P(s' \mid s_{1}, a) [\max_{a' \in \mathcal{A}} Q_{r,t+1(T)}^{*}(s', a')]$$

$$\stackrel{b}{=} r(s_{1}, a) + \gamma \sum_{z' \in \mathcal{Z}} \sum_{s' \in \Theta^{-1}(z')} P(s' \mid s_{1}, a) [\max_{a' \in \mathcal{A}} Q_{r,t+1(T)}^{*}(s', a')]$$

$$\stackrel{c}{=} r(s_{2}, a) + \gamma \sum_{z' \in \mathcal{Z}} \sum_{s' \in \Theta^{-1}(z')} P(s' \mid s_{2}, a) [\max_{a' \in \mathcal{A}} Q_{r,t+1(T)}^{*}(s', a')]$$

$$= Q_{r,t(T)}^{*}(s_{2}, a).$$
(1)

7)

850 851 852

Equalities (a) and (b) hold due to the surjective relationship between s and z. Equality (c) holds due 853 to the definition of bisimulation and the fact that, for $\forall z' \in \mathcal{Z}_{bisim}, \max_{a' \in \mathcal{A}} Q_{r,t+1}^*(s',a')$ is the 854 same for all $s' \in \Theta_{bisim}^{-1}(z')$, as established by the induction hypothesis. 855

856 We denote $Q_r^* = Q_{r,t(\infty)}^*$ as the optimal Q function for the infinite-horizon MDP. The uniqueness of Q_r^* is guaranteed by the fixed-point property of Bellman operator. For any $s \in S$, $a \in A$, and 858 timestep $t \in \{1, 2, ..., T-1\}$. The optimal Q value gap between Q_r^* and $Q_{r,t(T)}^*$ can be expressed as: 859

860

861
$$Q_r^*(s,a) - Q_{r,t(T)}^*(s,a)$$

 $Q_{r,t(T)}^{*}(s_1,a)$

 $-r(e, a) \pm \alpha \sum$

863
$$= \left| r(s,a) + \sum_{s' \in S} P(s'|s,a) [\gamma \max_{a' \in A} Q_r^*(s',a')] - \right|$$

$$= \left| \sum_{s' \in \mathcal{S}} P(s'|s, a) [\gamma \max_{a' \in \mathcal{A}} Q_r^*(s, a) - \gamma \max_{a' \in \mathcal{A}} Q_{r,t+1(T)}^*(s', a')] \right|$$

 $\overset{a}{\leq} \sum_{s' \in \mathcal{S}} P(s'|s, a) \big| \big[\gamma \max_{a' \in \mathcal{A}} Q_r^*(s', a') - \gamma \max_{a' \in \mathcal{A}} Q_{r, t+1(T)}^*(s', a') \big] \big|$ 870

8

we have:

where inequality (a) holds due to the triangle inequality property, and inequalities (b) and (c) 878 follow from the properties of the maximum function. Assuming that Q_r^* and $Q_{r,T(T)}^*$ are bounded 879 for any $s \in S$ and $a \in A$, we conclude that for a discount factor $\gamma \in (0,1)$, the difference 880 $|Q_r^*(s,a) - Q_{r,t(T)}^*(s,a)|$ converges to 0 as $T \to \infty$. Applying backward induction as introduced in Eq. 17, we deduce that, for any $s_1, s_2 \in S$ that are bisimilar (i.e., $\Theta_{bisim}(s_1) = \Theta_{bisim}(s_2)$), $Q_r^*(s_1, a) = Q_r^*(s_2, a)$. This completes the proof that $\Theta_{bisim} \succeq \Theta_{Q^*}$ for general Bellman operators. 883 For the safety Bellman operator, we have an analogous definition that for any timestep $t \in$ 884 $\{1, 2, ..., T\}, Q_{h,t(T)}^*$ is the optimal-Q function at timestep t. Then, for $\forall s \in S$ and $\forall a \in A$, 885

 $\overset{b}{\leq} \sum_{s' \in \mathcal{S}} P(s'|s, a) \gamma \max_{a' \in \mathcal{A}} \left| Q_r^*(s', a') - Q_{r, t+1(T)}^*(s', a') \right|$

(18)

 $\leq^{c} \gamma \max_{a' \in A} \sum_{s' \in S} |Q_{r}^{*}(s', a') - Q_{r,t+1(T)}^{*}(s', a')|,$

$$Q_{h,t(T)}^*(s,a) = (1-\gamma)h(s) + \gamma \sum_{s' \in \mathcal{S}} P(s' \mid s,a) [\max\{h(s), \min_{a' \in \mathcal{A}} Q_{h,t+1(T)}^*(s',a')\}].$$
(19)

888 889 890

891

892

893

894

886

887

We denote $Q_{h,T+1(T)}^*(s,a) = h(s)$ for all $s \in S$ and $a \in A$, which implies that $Q_{h,T(T)}^*(s,a) = h(s)$ if h(s) > 0 and $Q_{h,T(T)}^*(s,a) = (1 - \gamma)h(s)$ otherwise. Given any state $s \in S$, for all $z' \in \mathcal{Z}_{bisim}$, $\max\{h(s), \min_{a' \in \mathcal{A}} Q_{h,T(T)}^*(s', a')\}$ is the same for all $s' \in \Theta_{bisim}^{-1}(z')$. For any $s_1, s_2 \in \mathcal{S}$ that are bisimilar, and for any action $a \in A$, we apply backward induction from timestep T-1 to 1 such that:

895 896

$$Q_{h,t(T)}^{*}(s_{1},a) = (1-\gamma)h(s_{1}) + \gamma \sum_{s' \in \{\cup_{z' \in \mathcal{Z}} \Theta^{-1}(z')\}} P(s' \mid s_{1},a)[\max\{h(s_{1}), \min_{a' \in \mathcal{A}} Q_{h,t+1(T)}^{*}(s',a')\}] = (1-\gamma)h(s_{1}) + \gamma \sum_{z' \in \mathcal{Z}} \sum_{s' \in \Theta^{-1}(z')} P(s' \mid s_{1},a)[\max\{h(s_{1}), \min_{a' \in \mathcal{A}} Q_{h,t+1(T)}^{*}(s',a')\}] = (1-\gamma)h(s_{2}) + \gamma \sum_{z' \in \mathcal{Z}} \sum_{s' \in \Theta^{-1}(z')} P(s' \mid s_{2},a)[\max\{h(s_{2}), \min_{a' \in \mathcal{A}} Q_{h,t+1(T)}^{*}(s',a')\}] = Q_{h,t(T)}^{*}(s_{2},a).$$

$$(20)$$

Similarly, we denote $Q_h^* = Q_{h,t(\infty)}^*$ as the optimal Q function with safety Bellman operator for the infinite-horizon MDP. The uniqueness of Q_h^* is also guaranteed by the fixed-point property of safety 908 Bellman operator. For any $s \in S$, $a \in A$, and timestep $t \in \{1, 2, ..., T-1\}$. The optimal Q value 909 gap between Q_h^* and $Q_{h,t(T)}^*$ can be expressed as: 910

911

906

907

912 913

917

$$= \left| (1-\gamma)h(s) + \sum_{s' \in \mathcal{S}} P(s'|s,a) [\gamma \max\{h(s), \min_{a' \in \mathcal{A}} Q_h^*(s',a')\} - (1-\gamma)h(s) - \sum_{s' \in \mathcal{S}} P(s'|s,a) [\gamma \max\{h(s), \min_{a' \in \mathcal{A}} Q_{h,t+1(T)}^*(s',a')\} \right|$$

 $|Q_h^*(s,a) - Q_{h,t(T)}^*(s,a)|$

 $\leq \sum_{s' \in \mathcal{S}} P(s'|s, a) \gamma \max_{a' \in \mathcal{A}} \left| Q_h^*(s', a') - Q_{h, t+1(T)}^*(s', a') \right|$

 $\leq \gamma \max_{a' \in \mathcal{A}, s' \in \mathcal{S}} \big| Q_h^*(s', a') - Q_{h,t+1(T)}^*(s', a') \big|.$

918
919
$$= \left| \sum_{s' \in \mathcal{S}} P(s'|s, a) [\gamma \max\{h(s), \min_{a' \in \mathcal{A}} Q_h^*(s', a')\} - \gamma \max\{h(s), \min_{a' \in \mathcal{A}} Q_{h,t+1(T)}^*(s', a')\}] \right|$$
920

922

927 928

929

930

931

932

938 939

940 941 942

943

949

A similar conclusion can be drawn that for a discount factor $\gamma \in (0, 1)$, the difference $|Q_h^*(s, a) - Q_{h,t(T)}^*(s, a)|$ approaches to 0 as $T \to \infty$, provided that Q_r^* and $Q_{r,T(T)}^*$ are bounded for any $s \in S$ and $a \in A$. With backward induction as introduced in Eq. 20, we conclude that $\Theta_{bisim} \succeq \Theta_{Q^*}$ for safety Bellman operators.

 $\leq \sum_{s'\in\mathcal{S}} P(s'|s,a) \left| \gamma \max\{h(s), \min_{a'\in\mathcal{A}} Q_h^*(s',a')\} - \gamma \max\{h(s), \min_{a'\in\mathcal{A}} Q_{h,t+1(T)}^*(s',a')\} \right|$

Lemma A.1. For any MDP, given the optimal Q-functions induced by either the general Bellman operator or the safety Bellman operator, the optimal policy for Θ_{Q^*} remains optimal in the ground MDP.

Proof. For any state $s \in S$, we observe that $Q^*(\Theta_{Q^*}(s), a) = Q^*(s, a)$ for any $a \in A$, as per Definition 3.2. It is evident that for a given state $s \in S$:

$$a_r^* = \arg\max_{a \in \mathcal{A}} Q_r^*(\Theta_{Q_r^*}(s), a) = \arg\max_{a \in \mathcal{A}} Q_r^*(s, a), \tag{22}$$

(21)

$$a_h^* = \arg\min_{a \in \mathcal{A}} Q_h^*(\Theta_{Q_h^*}(s), a) = \arg\min_{a \in \mathcal{A}} Q_h^*(s, a).$$
⁽²³⁾

Therefore, we conclude that the optimal policy is preserved for Q^* -irrelevant representations.

Lemma A.2. For any MDP, given the optimal Q-functions induced by either the general Bellman operator or the safety Bellman operator, and any representation Θ_1 that is finner than Θ_{Q^*} , i.e. $\Theta_1 \succeq \Theta_{Q^*}$, it holds that $Q^*(\Theta_1(s), a) = Q^*(s, a)$ for any $s \in S$ and $a \in A$. The optimal policy for Θ_1 is also optimal in the ground MDP.

Proof. We denote the optimal value function for representation Θ_1, Θ_{Q^*} and ground state s as $Q_{\Theta_1}^*, Q_{\Theta_{Q^*}}^*$ and Q_s^* , respectively. It is evident that $Q_{\Theta_1}^*(\Theta_1(s), a) = Q_{\Theta_{Q^*}}^*(\Theta_{Q^*}(s), a) = Q_s^*(s, a)$ for any $s \in S, a \in A$ is one of the solutions for $Q_{\Theta_1}^*$, due to the subjective relationship between $\Theta_1(s)$ and $\Theta_{Q^*}(s)$. To show that the optimal value function for representation Θ_1 is unique, suppose that there exist two optimal value functions $Q_{r1}^*(z, a)$ and $Q_{r2}^*(z, a)$ for any representation $z \in Z_{\Theta_1}$ and action $a \in A$. The gap between them can be expressed as

$$\Delta_{Q_r^*}(z,a) = |Q_{r1}^*(z,a) - Q_{r2}^*(z,a)| = \left| \sum_{z'} P(z'|z,a) \gamma \left(\max_{a'} Q_{r1}^*(z',a') - \max_{a'} Q_{r2}^*(z',a') \right) \right|$$
(24)
$$\leq \gamma \max_{z',a'} \Delta_{Q_r^*}(z',a'),$$

where the inequality directly arises from the reasoning outlined in Eq. 18. With the discount factor $\gamma \in (0, 1)$, $\Delta_{Q_r^*}(z, a)$ tends to zero for any finite value of $Q_{r1}^*(z, a)$ and $Q_{r2}^*(z, a)$. This implies that the fixed point Q^* for the general Bellman operator is always unique.

For the safety Bellman operators, we also have

96 97 97

962

963

964

$$\Delta_{Q_{h}^{*}}(z,a) = |Q_{h1}^{*}(z,a) - Q_{h2}^{*}(z,a)|$$

$$= \left| \sum_{z'} P(z'|z,a) \gamma \left(\max\{h(z), \min_{a'} Q_{h1}^{*}(z',a')\} - \max\{h(z), \min_{a'} Q_{h1}^{*}(z',a')\} \right) \right|$$

$$\leq \gamma \max_{z',a'} \Delta_{Q_{h}^{*}}(z',a'),$$
(25)

where the inequality is a straightforward result of the proof sketch given in Eq. 21. It can be concluded that $Q^*(\Theta_1(s), a) = Q^*(s, a)$ holds for any $\Theta_1 \succeq \Theta_{Q^*}$, both for the general Bellman operator and the safety Bellman operator. Therefore, the optimal policy for Θ_1 is also optimal in the ground MDP, following the proof sketch provided in Lemma A.1.

Combining Eqs. 17 18 and 20 21, we conclude that $\Theta_{\text{bisim}} \succeq \Theta_{Q^*}$ holds for both the general Bellman operator and the safety Bellman operator. Combining Lemma A.1 and Lemma A.2, we conclude that both Θ_{bisim} and Θ_{Q^*} preserve optimality for the ground MDP. The proof of Theorem 3.1 is complete.

Theorem 3.1 shows that our Q^* -irrelavance representation leads to smaller representation space than bisimulation representation, while preserving the optimal policy. A smaller representation space typically implies higher generalization capabilities and higher sampling efficiency during the policy learning process.

Proposition A.3. For any MDP, $\Theta_1 \succeq \Theta_2$ indicates $H(s|\Theta_1(s)) \le H(s|\Theta_2(s))$.

Proof. We denote z_1 and z_2 as representations of $\Theta_1(s)$ and $\Theta_2(s)$, respectively. We have

$$H(s|z_{1}) = -\sum_{s,z_{1}} p(s, z_{1}) \log \frac{p(s, z_{1})}{p(z_{1})}$$

$$= -\sum_{s,z_{1}} p(s, z_{1}) \log p(s, z_{1}) + \sum_{s} \sum_{z_{1}} p(z_{1}) p(s|z_{1}) \log p(z_{1})$$

$$\stackrel{a}{=} -\sum_{s} p(s) \log p(s) + \sum_{z_{1}} p(z_{1}) \log p(z_{1})$$

$$\stackrel{b}{=} -\sum_{s} p(s) \log p(s) + \sum_{z_{2}} p(z_{2}) \sum_{z_{1}} p(z_{1}|z_{2}) \log p(z_{1})$$

$$\stackrel{c}{\leq} -\sum_{s} p(s) \log p(s) + \sum_{z_{2}} p(z_{2}) \log \sum_{z_{1}} p(z_{1}|z_{2}) p(z_{1})$$

$$\stackrel{d}{\leq} -\sum_{s} p(s) \log p(s) + \sum_{z_{2}} p(z_{2}) \log \sum_{z_{1}} \mathbb{I}_{\{p(z_{1}|z_{2})\neq 0\}} p(z_{1})$$

$$\stackrel{e}{=} -\sum_{s} p(s) \log p(s) + \sum_{z_{2}} p(z_{2}) \log p(z_{2})$$

$$= H(s|z_{2}), \qquad (26)$$

1005

1002

1004

986

987

where $\mathbb{I}_{\{p(z_1|z_2)\neq 0\}} = 1$ if $p(z_1|z_2) \neq 0$, and $\mathbb{I}_{\{p(z_1|z_2)\neq 0\}} = 0$ otherwise. Equality (a) holds as z_1 is a function of s. Note that $\sum_{z_1\in\mathcal{Z}_1} p(z_1|z_2) = 1$ for $\forall z_2\in\mathcal{Z}_2$ in equality (b). Inequality (c) is a consequence of Jensen's inequality. Inequality (d) holds since for $\forall z_2\in\mathcal{Z}_2$, the conditional probability $p(z_1|z_2)$ does not exceed 1 for all $z_1\in\{z\in\mathcal{Z}_1|p(z|z_2)\neq 0\}$. Equality (e) holds due to the surjective relationship between z_1 and z_2 .

Based on Theorem 3.1 and Proposition A.3, we conclude that:

1012 1013 1014

$$0 \le H(s|\Theta_{\text{bisim}}(s)) \le H(s|\Theta_{Q^*}(s)).$$
(27)

1015 Given our primary objective of maximizing the conditioned entropy $H(s|z_{\theta}(s))$, the proposed Q-1016 supervised contrastive learning method theoretically exhibits superior generalization capabilities 1017 compared to bisimulation.

1018 1019

1020

B METHODOLOGY CLARIFICATIONS

1021 B.1 IMPLICIT Q-LEARNING

Implicit Q-Learning (IQL) is the pioneering in-sample offline RL algorithm proposed by Kostrikov
 et al. (2021). It decouples the estimation of optimal Q-values from policy optimization, enabling
 implicit policy learning through the value function. Unlike standard Q-learning, which explicitly
 derives a policy by maximizing Q-values, IQL avoids direct maximization, reducing susceptibility

to instability issues such as overestimation or divergence. The core technique in IQL is expectile regression. Given a random variable X with an unknown distribution, the $\tau \in (0, 1)$ expectile can be estimated by solving:

$$\arg\min_{m_{\tau}} \mathbb{E}_{x \sim X}[L^{\tau}(x - m_{\tau})], \text{ where } L^{\tau}(u) = |\tau - \mathbb{I}(u < 0)|u^2$$
(28)

1032 Specifically, as $\tau \to 1$, the solution to Eq. 28 approximates the upper bound of the random variable 1033 X. Extending this to conditional distributions, the optimal value function can be approximated by 1034 minimizing:

$$\mathcal{L}_{V} = \mathbb{E}_{(s,a) \sim D_{\beta}} \left[L^{\tau}(Q(s,a) - V(s)) \right],$$
(29)

(30)

(33)

and the optimal Q function can be updated accordingly with the TD loss

1042

1035

1029 1030 1031

This process yields in-support optimal value and Q functions without training optimal policies. In
 this paper, we extend IQL with safety Bellman Operator and estimate both upper bound and lower
 bound of the value functions for safety assessments.

 $\mathcal{L}_Q = \mathbb{E}_{(s,a,s',r)\sim D_\beta} \left[(r + \gamma V(s') - Q(s,a))^2 \right].$

1043 1044 B.2 DIFFUSION BEHAVIOR CLONER

The diffusion model was initially introduced as an iterative denoising framework for image generation in the domain of *computer vision* (Sohl-Dickstein et al., 2015; Ho et al., 2020). More recently, it has been adapted for decision-making in state-based tasks, due to its superior performance in capturing action distributions within a dataset. As introduced in Section 3.2, SDQC requires a behavior cloner to capture and reproduce in-support actions for each state within the offline datasets. Following the approach of Lu et al. (2022a;b), we employ score-based diffusion and the DPM-Solver. The training loss for the behavior cloner π_{behav} is expressed as follows:

$$\mathcal{L}_{\pi_{behav}} = \mathbb{E}_{\mathbf{t} \sim \mathrm{U}(1,\mathrm{T}), \zeta \sim \mathcal{N}(0,\mathrm{I}), (\mathbf{s},\mathbf{a}) \sim \mathrm{D}_{\beta}}[\|\zeta - \zeta_{\psi_{behav}}(a_t,s,t)\|].$$
(31)

1054 After training converges, we use second-order DPM-Solver (Lu et al., 2022a) to form π_{behav} and 1055 sample $|\mathcal{A}_{\beta}^{s}|$ actions for each state in offline datasets \mathcal{D}_{β} . These actions will be utilized in the 1056 subsequent joint optimization of Q-functions and representations.

1057

1059

1063 1064

1067 1068 1069

1052 1053

1058 B.3 DIFFUSION POLICY

As described in Eqs. 13 and 14, we train three distinct diffusion policies using weighted regression (Zheng et al., 2024). In line with most existing diffusion-based policies (Wang et al., 2022; Garg et al., 2023; Lu et al., 2023), our three policies can be formulated as follows:

$$\begin{cases} \pi_{r}(a|z_{\theta_{r}}(s)) = p_{\psi_{r}}(a_{0:T}|z_{\theta_{r}}(s)) = \mathcal{N}(a_{T};\mathbf{0},\mathbf{I}) \prod_{t=1}^{T} p_{\psi_{r}}(a_{t-1}|a_{t},z_{\theta_{r}}(s)) \\ \pi_{h}(a|z_{\theta_{h}}(s)) = p_{\psi_{h}}(a_{0:T}|z_{\theta_{h}}(s)) = \mathcal{N}(a_{T};\mathbf{0},\mathbf{I}) \prod_{t=1}^{T} p_{\psi_{h}}(a_{t-1}|a_{t},z_{\theta_{h}}(s)) \\ \pi_{to}(a|z_{\theta_{r}}(s),z_{\theta_{h}}(s)) = p_{\psi_{to}}(a_{0:T}|z_{\theta_{r}}(s),z_{\theta_{h}}(s)) \\ = \mathcal{N}(a_{T};\mathbf{0},\mathbf{I}) \prod_{t=1}^{T} p_{\psi_{to}}(a_{t-1}|a_{t},z_{\theta_{r}}(s),z_{\theta_{h}}(s)), \end{cases}$$
(32)

where the reverse transitions are modeled as Gaussian process:

$$\begin{cases} p_{\psi_r}(a_{t-1}|a_t, z_{\theta_r}(s)) = \mathcal{N}(a_{t-1}; \mu_{\psi_r}(a_t, z_{\theta_r}(s), t), \Sigma(t)) \\ p_{\psi_h}(a_{t-1}|a_t, z_{\theta_h}(s)) = \mathcal{N}(a_{t-1}; \mu_{\psi_h}(a_t, z_{\theta_h}(s), t), \Sigma(t)) \\ p_{\psi_h}(a_{t-1}|a_t, z_{\theta_h}(s), z_{\theta_h}(s)) = \mathcal{N}(a_{t-1}; \mu_{\psi_h}(a_t, z_{\theta_h}(s), t), \Sigma(t)) \end{cases}$$

$$(p_{\psi_{to}}(a_{t-1}|a_t, z_{\theta_r}(s), z_{\theta_h}(s)) = \mathcal{N}(a_{t-1}; \mu_{\psi_{to}}(a_t, z_{\theta_r}(s), z_{\theta_h}(s), t), \Sigma(t))$$

1076 1077

Given a variance schedule defined by $\beta_t = 1 - \alpha_t$, we proceed to define:

1078
1079
$$\bar{\alpha}_t = \prod_{i=1}^t \alpha_i, \quad \tilde{\beta}_t = \frac{1 - \bar{\alpha}_{t-1}}{1 - \bar{\alpha}_t} \beta_t.$$
(34)

The mean of the Gaussian process is then given by:

1084

 $\begin{cases} \mu_{\psi_r}(a_t, z_{\theta_r}(s), t) = \frac{1}{\sqrt{\alpha_t}} (a_t - \frac{\beta_t}{\sqrt{1 - \bar{\alpha}_t}} \zeta_{\theta_r}(a_t, z_{\theta_r}(s), t)) \\ \mu_{\psi_h}(a_t, z_{\theta_h}(s), t) = \frac{1}{\sqrt{\alpha_t}} (a_t - \frac{\beta_t}{\sqrt{1 - \bar{\alpha}_t}} \zeta_{\theta_h}(a_t, z_{\theta_h}(s), t)) \\ \mu_{\psi_{to}}(a_t, z_{\theta_r}(s), z_{\theta_h}(s), t) = \frac{1}{\sqrt{\alpha_t}} (a_t - \frac{\beta_t}{\sqrt{1 - \bar{\alpha}_t}} \zeta_{\theta_{to}}(a_t, z_{\theta_r}(s), z_{\theta_h}(s), t)), \end{cases}$

(35)

and the covariance matrix is expressed as
$$\Sigma(t) = \tilde{\beta}_t I$$
. During the testing phase, actions can be sampled from the reverse diffusion chain for each diffusion policy. To account for safety considerations, we sample multiple actions and select the one with the lowest Q_h value as the final action to be executed (cf. Appendix E.2 for details).

1091

1093

C IMPLEMENTATION DETAILS

C.1 SDQC SUMMARIZATION

To provide an intuitive understanding the mechanism of SDQC, we present a brief summary of its training and deployment scheme in this subsection.

1099 Algorithm 1 SDQC Training

1100 Phase 1: Behavior cloner training and behavior actions generation **Require:** Initial network $\zeta_{\psi_{behav}}$, datasets $\mathcal{D}_{\beta} = \{s, a, s', r, h\}^N$ 1101 1: for each iteration $\{s, a\}^{M_1} \sim \mathcal{D}_\beta$ do 1102 2: Update $\zeta_{\psi_{behav}}$ using Eq. 31 1103 3: end for 1104 4: for each state $s \sim D_{\beta}$ do 1105 5: Given s, generate multiple behavior actions \mathcal{A}^s_{β} with $\zeta_{\psi_{bchav}}$ using DPM-Solver (Lu et al., 2022a) 1106 6: end for 7: **return** Updated datasets $\mathcal{D}_{\beta} = \{s, \mathcal{A}_{\beta}^{s}, a, s', r, h\}^{N}$ 1107 Phase 2: Joint optimization for value functions and representations 1108 **Require:** Initial reward-related network V_r, Q_r, z_{θ_r} , cost-related network $V_h^{\text{up}}, V_h^{\text{low}}, Q_h, z_{\theta_h}$, and datasets 1109 $\mathcal{D}_{\beta} = \{s, \mathcal{A}^{s}_{\beta}, a, s', r, h\}^{N}$ 1110 8: for each iteration $\{s, \mathcal{A}_{\beta}^{s}, a, s', r\}^{M_{2}} \sim \mathcal{D}_{\beta}$ do 9: Update $V_{r}, Q_{r}, z_{\theta_{r}}$ jointly using Eq. 8 1111 1112 10: end for 11: for each iteration $\{s, \mathcal{A}^s_{\beta}, a, s', h\}^{M_2} \sim \mathcal{D}_{\beta}$ do 1113 1114 Update $V_h^{\text{up}}, V_h^{\text{low}}, Q_h, z_{\theta_h}$ jointly using Eq. 12 12: 1115 13: end for 14: return Reward/Cost-related value function and representation networks $V_r, Q_r, z_{\theta_r}, V_{\mu}^{\rm up}, V_h^{\rm low}, Q_h, z_{\theta_h}$ 1116 **Phase 3: Three policies extraction** 1117 **Require:** Initial policy network $\zeta_{\psi_r}, \zeta_{\psi_h}, \zeta_{\psi_{to}}$, fixed pre-trained network $V_r, Q_r, z_{\theta_r}, V_h^{\text{low}}, Q_h, z_{\theta_h}$, and 1118 datasets $\mathcal{D}_{\beta} = \{s, a\}^N$ 1119 15: for each iteration $\{s, a\}^{M_3} \sim \mathcal{D}_\beta$ do 1120 Calculate regression weight for three policies with $V_r, Q_r, z_{\theta_r}, V_h^{\text{low}}, Q_h, z_{\theta_h}$ using Eq. 14 16: 1121 17: Update $\zeta_{\psi_r}, \zeta_{\psi_h}, \zeta_{\psi_{to}}$ using Eq. 13 1122 18: end for 19: **return** Three distinct policies π_r, π_h, π_{to} 1123 1124 1125 As introduced in Section 3, SDQC requires a three-phase training process (see Algorithm 1). The 1126 primary objective of the first stage (lines 1-7) is to generate a set of behavior actions \mathcal{A}^s_β for each 1127 state $s \in \mathcal{D}_{\beta}$. The set is then utilized in the second stage to measure the similarity between states, represented by $d(s, \tilde{s}) \approx \sup_{a \in \mathcal{A}_{\beta}^{s}} |Q^{*}(z_{\theta}(s), a) - Q^{*}(z_{\theta}(\tilde{s}), a)|$. The second phase (lines 8-14), 1128 1129 known as the joint optimization of value functions and representations, employs expectile regression 1130 (Kostrikov et al., 2021) to learn in-support optimal value functions while simultaneously using contrastive learning to cluster similar states (measured by $d(s, \tilde{s})$) within the representation space. 1131 The final phase (lines 15-19) introduces three distinct policies, which are trained using weighted 1132 regression diffusion models as introduced by Zheng et al. (2024). It is important to note that these 1133

policies are conditioned on the representation space rather than the ground-truth state space. During

deployment, the selection of the specific policy to be adopted is guided by the evaluation of statesusing the cost-related value function. For further details, please refer to the subsequent paragraph.

1136 1137

1138

Algorithm 2 SDQC Deployment

Require: Three policies π_r, π_h, π_{to} , representation network $z_{\theta_h}, z_{\theta_r}$, and cost-related value-functions $V_h^{\text{up}}, V_h^{\text{low}}, Q_h$ 1140 1: **Given** any state *s* 1141 2: **if** $V_h^{\text{up}}(z_{\theta_h}(s)) > 0$ **then** 1142 (There exist in-support actions that may lead to unsafe outcomes) 1143 3: **if** $V_h^{\text{down}}(z_{\theta_h}(s)) > 0$ **then**

1143 3: If $V_h^{\text{barr}}(z_{\theta_h}(s)) > 0$ then 1144 4: sample actions: $\{a_i\}^{\text{cand.}} \sim \pi_h(\cdot | z_{\theta_h}(s))$

11444:sample actions5:else

1145 6: sample actions: $\{a_i\}^{\text{cand.}} \sim \pi_{to}(\cdot | z_{\theta_h}(s), z_{\theta_r}(s))$

7: end if

8: else (The state is absolute safe)

9: sample actions: $\{a_i\}^{\text{cand.}} \sim \pi_r(\cdot | z_{\theta_r}(s))$

1149 9: sai 10: end if

1150 11: return Final action $\arg \min_{a \in \{a_i\}^{\text{cand.}}} Q_h(z_{\theta_h}(s), a)$

1151 1152

1146

1147

1148

Upon completion of training, SDQC is deployed with three distinct policies π_r, π_h, π_{to} derived 1153 from the third training phase, along with two representation networks $z_{\theta_h}, z_{\theta_r}$ and three cost-related 1154 value functions $V_h^{\text{up}}, V_h^{\text{low}}, Q_h$ from the second training phase (cf. Algorithm 2). Given any state 1155 s, we initially conduct safety assessments based on cost-related representations. The condition 1156 $V_h^{\rm up}(z_{\theta_h}(s)) > 0$ (line 2) suggests the existence of in-support actions that might lead to unsafe 1157 outcomes, necessitating a joint consideration of safety and reward. A positive $V_h^{\text{down}}(z_{\theta_h}(s))$ (line 1158 3) indicates that no action can ensure safety in future trajectories; the agent's primary objective is 1159 therefore to exit the unsafe region by deploying policy π_h , regardless of the reward considerations. 1160 Conversely, $V_h^{\text{low}} \le 0 < V_h^{\text{up}}$ (line 5) reflects a borderline safe condition, requiring the agent to consider both reward and cost, thereby deploying policy π_{to} . On the other hand, $V_h^{\text{up}}(z_{\theta_h}(s)) \le 0$ 1161 1162 (line 8) confirms absolute safety, obviating the need for the agent to consider cost-related information. 1163 An illustration diagram is presented in the right subplot of Figure 2. Note that single action sampled 1164 by diffusion model is not trustworthy enough, thereby sampling batch actions and conduct the one with the lowest $Q_h(z_{\theta_h}(s), a)$ leads to safer outcomes (line 11). Please refer to Appendix E.2 for 1165 further details. 1166

1167

1168 C.2 SDQC NETWORK STRUCTURE

1169 As described in Section 3.3, the representa-1170 tions in our proposed SDQC framework are 1171 trained concurrently with the optimal Q value 1172 learning process. The neural network structure 1173 illustrated in Figure 5 is utilized for training 1174 both the reward- and cost-related representa-1175 tions as well as the value functions. The global 1176 observation s is encoded into the representa-1177 tion z, and the value functions (both V and 1178 Q) are computed based on this representation 1179 with separate multiple-layer-perceptron (MLP) neural networks. 1180

In certain safe RL benchmark problems, it is
observed that the majority of dimensions in the
global observation share similar physical meanings. For instance, in the Safety-Gymnasium
domain, a significant number of dimensions in

 $S \rightarrow \underbrace{\text{State Encoder}}_{(\text{MLP or ATN)}} + \underbrace{z}_{\text{MLP}} + \underbrace{z}_{\text{MLP}} + \underbrace{v}_{\text{MLP}} +$

Figure 5: Neural network structure for training value functions and representations of SDQC.

the global observations correspond to lidar measurements, which provide information about the distances between the agent and the destination or obstacles in specific directions. This reminds us of the self-attention mechanism (Vaswani et al., 2017), which is known for its superior ability to capture

relationships among input information that share similar representations in comparison to traditional
 MLP architectures.

Nevertheless, attention mechanisms typically rely on vector multiplication to compute attention weights, which presents a challenge when dealing with global observations where each dimension contains scalar information. Towards this end, we propose to transform each scalar observation dimension into a vector representation using a fixed Gaussian Fourier Encoder (Ho et al., 2020).
Subsequently, attention is applied to the encoded vector representations. The output of the attention module is then flattened and passed through an MLP to obtain the final representation. Please refer to Figure 6 for a detailed illustration of the network structure.



Figure 6: Neural network structure of the attention-based state encoder.

To demonstrate the efficacy of our attention-based state encoder (ATN) in effectively capturing infor-1210 mation from global observations and identifying the relevance of specific dimensions to reward/cost, 1211 we present the attention patterns (i.e., softmax $(Q \cdot K^T / \sqrt{d_k}))$ of the reward- and cost-related state 1212 encoders in the task "PointGoal2" (cf. Figure 7). In "PointGoal2", the global observations consist 1213 of 60 dimensions. Among them, the first 12 dimensions represent the self-status of the agent, the 1214 subsequent 16 dimensions contain reward-related information, and the last 32 dimensions contain 1215 cost-related information. Ideally, a reward-related attention pattern should assign higher attention 1216 weights to the first 28 dimensions while ignoring the last 32 dimensions. On the other hand, a 1217 cost-related attention pattern should focus on the first 12 and last 32 dimensions while disregarding 1218 the middle 16 dimensions. The observed attention patterns during our experiments align with the 1219 relevance of specific dimensions to reward and cost. For the ablation study on the state encoder network structure, please refer to Appendix D. 1220



Figure 7: Attention pattern of the reward- and cost-related state encoder, and the meaning of each observation dimension in the task "PointGoal2" from Safety-Gymnasium. Darker colors represent higher values. The "Goal Lidar" dimensions observe the position of the destination, indicating reward information. The "Hazard Lidar" and "Vases Lidar" dimensions observe the position of obstacles, indicating cost information. The pattern is averaged over 3000 observations randomly chosen from DSRL datasets.

1240

1208

1209

1241 C.3 SDQC HYPERPARAMETERS

1 4- 7 4-	As discussed in Section 3.2, to calculate the soft similarity
1243	measure for contrastive learning, our SDQC framework
1244	requires pre-training a generative model to capture the
1245	behavior policy of the offline datasets π_{β} . For this purpose,
1246	we employ diffusion probabilistic models (DPM) (Ho et al.,
1247	2020; Song et al., 2020) and utilize the DPM-Solver, a fast
1248	high-fidelity ODE solver proposed by Lu et al. (2022a;b),
1249	to generate the behavior actions of each state s in the offline
1250	datasets, denoted as \mathcal{A}^s_{β} . We utilize the default network
1251	configurations outlined in (Lu et al., 2022a; 2023). Specific
1050	hyperparameter settings can be found in Table 2.

Table 2: Hyperparameters of the DPM-solver for generating behavior actions.

Hyperparameters	Value
Learning rate	3e-4
Batch size	4096
Training steps	5e5
Diffusion timesteps	15
Generated action numbers	0
for each state $ \mathcal{A}_{\beta}^{s} $	δ

SDQC can simultaneously train the reward and cost valuefunctions and their respective representations using Eqs. 8

and 12. The network structures are described in detail in Appendix C.2, and generic hyperparameters can be found in Table 3. Regarding the updating of the safety Bellman operator, we follow the settings in FISOR (Zheng et al., 2024). The constraint violation function is defined as h(s) = -1 when the cost function satisfies c(s) = 0, and h(s) = 25 when c(s) > 0.

1259

10/10

1260 1261 1262

Table 3: Generic hyperparameters of SDQC in value functions and representations training phase.

1263	Module		Hyper-parameters	Value
1264			Optimizer	Adam
1265			Learning rate	3e-4
1266			Batch size	512
1267	General		Training steps	5e5
1268			Soft measure temperature factor η	1.0
1269			Contrastive temperature factor ν	0.1
1270			Contrastive term coef δ	1.0
1071			Number of hidden layers (Q & V)	2
1271			Number of neurons in hidden layer (Q & V)	256
1272	Critic		Activation function (Q & V)	Mish
1273	Chuc		Expectile $ au$	0.9
1274			Discount factor γ	0.99
1275			Target critic soft update	0.005
1276			Number of hidden layers	2
1277		MLP	Number of neurons in hidden layer	256
1278	State Encoder		Activation function	Mish
1279	State Liteoder		Number of head	2
1280		ATN	Embed dimension for each head	64
1281			Dropout rate	0.1
· · · · · · ·				

1282 1283

Considering the significant variation in physical meanings among the observation dimensions of different tasks, we employ different state encoder structures accordingly. For tasks that have observation dimensions with diverse physical meanings, we utilize the MLP structure. Conversely, for tasks
 where most observation dimensions have consistent physical meanings, we utilize the attention-based state encoder (ATN).

It is observed that the performance of SDQC with an ATN-based state encoder improves when trained with a larger contrastive loss coefficient (δ in Eqs. 8 and 12) and a higher number of anchor points ($|\mathcal{I}|$ in Eq. 5). On the other hand, the SDQC performs better with smaller values of δ and fewer anchor points if the MLP based encoder is used. Besides, the global observation dimensions vary across different tasks. For the ATN-based state encoder, we select the encoded state dimension (i.e., the dimensionality of z) to be approximately half of the global observations. On the other hand, for MLP, we choose the encoded state dimension to be roughly twice the size of the global observations. Optimal hyperparameters achieving the best performance on different tasks are presented in Table 4.

Domain	Agent	Task	State Encoder	Encoded State Dim	Contrastive Loss Coef	Anchor Number
Sofaty	Point Car	All All	ATN	32	1.0	8
Gymnasium	HalfCheetah Ant	Vel Vel	MLP	32	0.5	4
	Swimmer	Vel				
	Ball	All		16		
Bullet Safety	Car All Drone All		MLP	32 64	0.1	4
Survey	Ant	All		64		

Table 4: Hyperparameters of SDQC for different tasks. "All" denotes all different tasks for the same agent, while "Vel" refers to the velocity task.

For the final training phase, which involves policy extraction using weighted regressed diffusion models as described in Eq. 13, we follow the network structure design and generic diffusion parameter selection described by Zheng et al. (Zheng et al., 2024). We train three separate policies ($\pi_r \pi_h$ and π_{to}) with a learning rate of 0.0003, a batch size of 1024, and the total number of training steps is set to 500,000. The temperature parameters that control the strength of behavior regularization (in Eq. 14) are chosen as $\iota_r = \iota_{to} = 3.0$ and $\iota_h = 5.0$.

1317

1310

1296

1318 D ADDITIONAL ABLATION STUDIES

1319 Ablation studies on network structure. To demonstrate the effectiveness of our proposed attention-1320 based state encoder (cf. Appendix C.2), we conduct ablation studies on the neural network architec-1321 ture, as depicted in Figure 8. By substituting the attention-based state encoder with an MLP-based 1322 counterpart, we observe a deterioration in the performance of the SDQC, in terms of diminished 1323 rewards and increased costs. The t-SNE visualization results depicted in Figure 8b demonstrate 1324 that while the MLP-based state encoder does cluster representations with similar values in the 1325 high-dimensional space, the clustering effect is not as robust as that achieved by the attention-based 1326 approach. Consequently, this leads to an overestimation of the cost value, resulting in inaccurate 1327 assessments of the safety condition.



Figure 8: Ablation studies on the network structure in CarGoal2. (a) The actor-training-process evaluations of SDQC with attention-based (ATN) and MLP-based (MLP) state encoder. (b) t-SNE visualization of the distribution of the original state, the reward- and cost-related representations with ATN/MLP state encoders across 10 different safe trajectories. Anchor

Number

16

8

4

1

0

1350 Ablation studies on anchor number choice. An essential hyperparameter in our proposed Q-1351 supervised contrastive learning method is the anchor number, $|\mathcal{I}|$ in Eq. 5. This parameter determines 1352 the number of representation pairs to be clustered in the high-dimensional space during each gradient 1353 step. The ablation study results are summarized in Table 5. Our experimentation reveals that an 1354 anchor number can result in subpar clustering outcomes, consequently impairing the performance of SDQC. Conversely, overly large anchor numbers lead to an excessive influence of the contrastive loss 1355 term in the overall loss function, increasing the computational costs. To strike a balance and attain 1356 optimal performance, we let $|\mathcal{I}| = 8$ for our attention-based state encoders. 1357

CarGoal2

Cost

0.00

0.00

0.05

0.13

0.86

Reward

0.22

0.23

0.20

0.15

0.05

Table 5: Ablation studies on the choice of anchor number.

Runtime

(s/epoch)

31.4

28.8

26.0

23.6

19.4

CarPush2

Cost

0.06

0.04

0.06

0.18

2.15

Reward

0.28

0.31

0.25

0.10

0.21

Runtime

(s/epoch)

36.2

32.7

29.9

27.6

23.7

1358 1359

1361 1362 1363

1360

- 1364 1365
- 1366
- 1367 1368
- 1369

Ablation studies on contrastive-related hyperparameters. In our SDQC framework, one of the 1370 critical components is the contrastive representation loss, as described in Equations 8 to 12. This 1371 involves selecting appropriate values for the term coefficient δ and the exponential temperature ν . 1372 As shown in Tables 3 and 4, we vary δ across different domains but maintain a consistent $\nu = 0.1$ 1373 across all environments. The effects of these parameter choices are detailed in Table 6. With respect 1374 to the temperature ν , employing a very small value (0.01) tends to destabilize the training process, 1375 ultimately resulting in collapse. Conversely, using a larger value (1.0) produces poorly clustered 1376 representations, leading to a marked degradation in performance. Regarding the term coefficient δ , 1377 a smaller value results in a slight performance decline. However, a larger coefficient excessively 1378 prioritizes the contrastive loss, destabilizing the training of the value function and significantly degrading performance. While fine-tuning these hyperparameters for specific environments and tasks 1379 could potentially yield better experimental results on the benchmark, we choose not to do so. 1380

1381 1382

Table 6: Ablation studies on contrastive-related hyperparameters.

Env	Contrast Coef. (δ)	Contrast Temp. (ν)	Reward	Cost	Env	Contrast Coef.	Contrast Temp.	Reward	Cost
	1	0.01	NaN	NaN		1	0.01	NaN	NaN
	1	0.1	0.31	0.04		1	0.1	0.29	0.09
PointGoal2	1	1	0.22	0.16	CarPush2	1	1	0.20	0.48
	0.1	0.1	0.31	0.17		0.1	0.1	0.30	0.10
	10	0.1	0.23	0.48		10	0.1	0.22	0.53

1392 Ablation studies on the deployment of three distinct policies As introduced in Section 3.3 and 1393 Appendix C.1, SDQC coordinates three distinct policies, reward policy π_r , trade-off policy π_{to} , and 1394 cost policy π_h , to ensure excellent safety performance. To verify the necessity of each policy, we 1395 conduct ablation studies examining their individual deployments, with results presented in Table 7. 1396 Notably, a naive reward policy π_r focuses solely on maximizing rewards while ignoring costs, a naive cost policy π_h prioritizes minimizing costs but disregards rewards, and a naive trade-off policy π_{to} 1398 takes both into account but fails to excel in either maximizing rewards or minimizing costs. The 1399 best performance consistently results from the collaboration of all three policies. When the trade-off 1400 policy π_{to} is omitted (combining π_r and π_h), the agent incurs higher costs as it cannot respond 1401 promptly to borderline dangers. Combining the trade-off policy π_{to} and cost policy π_h does not increase costs, but results in a decline in reward accumulation. While combining the reward policy 1402 π_r and trade-off policy π_{to} achieves comparable performance to using all three policies, it results in 1403 slightly higher costs due to the agent's reduced ability to quickly escape dangerous situations.

	-	v	2
а.	л	\sim	1

÷	л	n	7
1	4	V	1

1408 1409

1410

- 1411
- 1412 1413

1414

Table 7: Ablation studies on the deployment of three distinct policies.

Policy Num. One							Two						Three		
Env Name	Naïve	Naïve	e π_{to} Naïve 7		π_h	π_h π_r and π_{to}		π_r and π_h		π_{to} and π_h		$\pi_r, \pi_{to} \text{ and } \pi_h$			
Env Maine	reward	cost	reward	cost	reward	cost	reward	cost	reward	cost	reward	cost	reward	cost	
PointGoal1	0.69	4.92	0.27	0.69	0.01	0.00	0.32	0.42	0.34	0.51	0.22	0.18	0.35	0.36	
PointGoal2	0.75	13.28	0.28	0.18	-0.11	0.00	0.23	0.24	0.20	0.14	0.23	0.12	0.29	0.09	
CarPush1	0.38	1.32	0.26	0.00	0.05	0.00	0.27	0.00	0.27	0.12	0.21	0.00	0.30	0.00	
CarPush2	0.42	4.34	0.27	0.01	0.02	0.00	0.31	0.01	0.28	0.16	0.18	0.03	0.31	0.04	

E EXPERIMENTAL DETAILS

1415 E.1 TASK DESCRIPTION

1417 **Safety-Gymnasium (Ray et al., 2019).** A collection of environments based on the Mujoco physics 1418 simulator. In the obstacle avoidance series environments, there are two agents (Point and Car) and 1419 three main tasks (Goal, Button, and Push), each with two levels of difficulty (1 and 2). Agents aim 1420 to reach the goal while avoiding any contact with obstacles. The environments are named using 1421 the following convention: {Agent}{Task}{Difficulty}. In the velocity-constrained environments, 1422 there are three agents: Ant, HalfCheetah, and Swimmer. The primary objective of these agents is to maximize their rewards while adhering to the imposed velocity constraints. The environments are 1423 named in the convention of {Agent} Velocity. 1424

1425

Bullet-Safety-Gym (Gronauer, 2022). A suite of environments built upon the PyBullet physics
simulator. These environments are similar to Safety-Gymnasium but feature a broader range of agents
(including Ball, Car, Drone, and Ant). The tasks are relatively straightforward, with only two options
available (Circle and Run). The environments are named through {Agent}{Task}.

1430

1431 E.2 EXPERIMENT SETTINGS 1432

We train the baseline algorithms using the recommended hyperparameters specific to each task, for BCQ-Lag (Fujimoto et al., 2019; Stooke et al., 2020), CPQ (Xu et al., 2022a), COptiDICE (Lee et al., 2022), CDT (Liu et al., 2023b), TREBI (Lin et al., 2023), and FISOR (Zheng et al., 2024). To ensure a fair comparison, we train the baseline algorithms with three different random seeds and save the final output policy for safety evaluation. For each output policy, we conduct evaluations over 20 episodes to obtain reliable performance measures.

As for the training process of SDQC, we follow the neural network structure design and hyperparameter settings in Appendix C. Analogously, we select three different random seeds for training and perform evaluation over 20 episodes for each seed. To improve safety performance, we follow Zheng et al. (2024) to sample 16 candidate actions for each RL timestep, regardless of the safety assessment results and policy usage. The safest action is then selected based on the lowest $Q_h(z_{\theta_h}(s), a)$ value and executed as the final action.

- 1445
- 1446 E.3 COMPUTATIONAL COSTS
- 1447

1448 We implement SDQC using PyTorch (Paszke et al., 2019) and conduct experiments on a single machine equipped with one GPU (NVIDIA RTX 4090, 24GB) and one CPU (AMD Ryzen 9 1449 7950X). The training process comprises three phases. The first phase, known as the diffusion 1450 behavior cloner, demands approximately 1 hour for each task. For the second phase, which involves 1451 training representations and critics, the duration varies depending on the chosen network architecture. 1452 Attention-based architectures typically require over 4 hours, whereas MLP-based architectures 1453 typically demand around 1 hour. Finally, in the last training phase, the diffusion actor, convergence 1454 typically occurs in about 1 hour (without online testing). 1455

Besides, we assess the inference time consumption of all baseline algorithms across 1000 RL
 timesteps on the CarPush2 task, averaging the results over 10 trials. Although SDQC is relatively slower compared to other non-autoregressive policies, it remains within an acceptable range.

Table 8: Inference time (seconds) comparison over 1000 RL timesteps.

BCQ-Lag	CPQ	COptiDICE	CDT	TREBI	FISOR	SDQC
1.51	1.85	1.86	3.45	585.87	6.11	11.13

1461 1462 1463

1464

1465

1499 1500

1458

1459 1460

F ADDITIONAL EXPERIMENTAL RESULTS

1466 F.1 Additional Generalization tests 1467 F.1 Additional Generalization tests

In addition to the generalization tests (on the agent "Car") presented in Section 4.2, we perform generalization tests on the "PointGoal" and "PointPush" tasks (in Safety-Gymnasium), as illustrated in Figure 9. Similarly, the "Point" agent is challenged with tasks that involve reaching a goal point or pushing a box to a goal point in hazardous areas with obstacles, with the difficulty level varying between simple (PointGoal1, PointPush1) and challenging (PointGoal2, PointPush2).

1473 Experimental observations reveal that, in comparison to the "Car" agent, the "Point" agent demon-1474 strates a higher degree of inertia during its motion within the environment. Specifically, the "Point" 1475 agent lacks the ability to instantaneously halt or promptly alter its direction, thereby rendering the maintenance of safety more challenging in equivalent tasks when compared to the "Car" agent. This 1476 significantly undermines the generalization capability of most algorithms on the "Point" agent. For 1477 instance, the state-of-the-art (SOTA) safe offline RL algorithm FISOR performs poorly on the "Point" 1478 agent, exhibiting high costs across multiple environments. In contrast, our SDQC algorithm still 1479 achieves nearly zero violations in the majority of environments. 1480



Figure 9: The generalization tests on the agent "Point" in Safety-Gymnasium.

1501 1502 F.2 IMPACT OF REPRESENTATION LOSS ON VALUE ESTIMATIONS

A major concern for SDQC is that the joint optimization process of the value functions and represen-1503 tations (refer to Eqs. 8 and 12) may lead to instability in value estimation, subsequently affecting final 1504 performance. We analyze the performance of SDQC, both with and without representation loss, and 1505 compare it to FISOR by examining their respective Critic (Q) loss and Value loss patterns in relation 1506 to reward/cost metrics. This comparative analysis is conducted throughout the training process using 1507 two tasks: 'CarPush2' and 'BallCircle', with results illustrated in Figure 10. The experimental results indicate that while the inclusion of representation loss does lead to an increase in critic and 1509 value loss, it does not compromise the overall stability of the training. Furthermore, our proposed 1510 neural network architecture (used for "CarPush2"), with the incorporation of an attention-based state 1511 encoder, markedly improves the precision and stability of value function learning compared to the simple MLP utilized by FISOR.



Figure 10: Comparison of value function loss during the training process among SDQC with and without representation loss, and FISOR, on the 'CarPush2' and 'BallCircle' tasks.

1540 F.3 WHY NOT BISIMULATIONS

1541 As discussed in Section 3.4, the Q-supervised representation learning approach in SDQC is theoreti-1542 cally superior to bisimulation by introducing a coarser representation space. From an experimental 1543 perspective, we aim to verify this assertion. However, challenges arise during the initial training 1544 phase for bisimulation. Bisimulation typically involves an additional step of training a model-based 1545 estimator to learn state transition and reward/cost functions. We find that estimating cost functions is particularly challenging due to their non-smooth and sparsely distributed nature. This issue is espe-1546 cially pronounced in the Safety-Gymnasium domain. To investigate this, we conduct the following 1547 experiments. 1548

We treat state-action pairs with non-zero costs as positive samples and use binary cross-entropy loss for training. The estimation is based on representations enhanced by an additional 2-layer MLP, with 256 hidden neurons. The batch size is 256, and the learning rate is set at 3e-4. Training is conducted over 300,000 steps. Upon completion, we evaluate the estimation results as shown in Table 9. There are four possible outcomes in the relationship between the estimation and the ground-truth labels for each state-action pair: "True Positive (TP)", "False Positive (FP)", "True Negative (TN)", and "False Negative (FN)". We report the following metric for the final cost-function approximation:

Accuracy : $\frac{TP + TN}{TP + FP + TN + FN}$

1560

$$\operatorname{Recall}: \frac{\operatorname{TP}}{\operatorname{TP} + \operatorname{FN}} \qquad \operatorname{F1Score}: \frac{2 \times \operatorname{Precision} \times \operatorname{Recall}}{\operatorname{Precision} + \operatorname{Recall}}$$

 $Precision: \frac{TP}{TP + FP}$

While the model achieves a high overall prediction accuracy of 98%, this metric is potentially misleading due to significant class imbalance in the dataset, where positive samples represent only 7% of the total observations. Our primary focus is on samples with non-zero cost values, i.e., positive samples. The "precision" metric, which indicates the proportion of correctly identified positive predictions among all positive predictions made by the model, reaches approximately 85%. "Recall" indicates the probability that actual positive samples are correctly identified by the model, which

is about 77%, suggesting that the model's ability to assess dangerous conditions is insufficient to support subsequent bisimulation training. Notably, unlike bisimulation, our SDQC framework learns representations based on non-sparse and continuous optimal Q-functions, which are significantly easier to optimize compared to sparse and non-smooth cost functions, thereby implicitly addressing this issue.

Table 9: Cost function approximation for "PointGoal2" and "CarPush2" Tasks. States with non-zero cost values are identified as positive samples. We report the accuracy (Acc.), precision (Pre.), recall, and F1 score of the estimation results. To align with the bisimulation training approach, we report the estimation results for the cost function when trained without (w/o) and with (w) the transition model.

	Emr	Dataset	Pos.	w/o Transition Estimation				w Transition Estimation			
	EIIV		Ratio	Acc.	Pre.	Recall	F1	Acc.	Pre.	Recall	F1
	PointGoal2	Train	7.65%	0.980	0.859	0.782	0.819	0.979	0.862	0.770	0.814
		Test	7.74%	0.976	0.838	0.760	0.797	0.976	0.844	0.750	0.794
	CarPush2	Train	7.25%	0.977	0.890	0.784	0.834	0.976	0.862	0.795	0.827
		Test	7.33%	0.974	0.874	0.758	0.812	0.973	0.843	0.774	0.807







As complementary to the SOTA safe offline RL algorithm FISOR (Zheng et al., 2024), our SDQC
employs the same implicit Q-learning method (Kostrikov et al., 2021) to learn optimal value functions
and utilizes the safety Bellman operator in (Fisac et al., 2019) for safety assessment. Additionally,
we adopt their approach for policy extraction through training a weighted regressed diffusion model.
However, it should be noted that our decision-making process is based on decoupled representations
rather than global observations. Furthermore, our policies are completely decoupled, and different
policies are selected based on varying safety assessment results. To provide futher comparisons

1651

1652 1653

1623 1624 Reward (SDQC) ---- Reward (FISOR) Cost (SDQC) ---- Cost (FISOR) ----1625 Safety-Gymnasium-PointGoal1 Safety-Gymnasium-PointGoal2 1626 0.35 0.8 0.30 0.7 3.5 3.5 1627 6. 0.6 st 0 3.0 2.5 2.0 2.0 P 0.6 0.25 Reward 0.20 1628 g 0.5 2.5 0.15 alized R 0.3 alized 1629 2.0 0.10 1.5 E 1.5 0.05 1630 U 0.2 N 1.0 ອັ້ 1.0 Norn 0.00 1631 0.1 0.5 0.5 -0.05 0.0 1632 0.0 -0.10 0.0 100 200 300 400 100 200 300 400 ò 500 ò 100 200 300 400 500 ò 100 200 300 400 500 ò 500 adient steps (×10³) adient steps (×10³) nt steps (×10³) Gradient steps (×10³) 1633 Safety-Gymnasium-PointPush1 Safety-Gymnasium-PointPush2 0.35 0.35 4.0 4 0 1634 0.30 0.30 3.5 3.5 1635 0.6 st 0.25 3.0 2.5 Reward 0.25 0.20 Rev 0.20 2.5 2.5 2.0 1.5 0.15 2.0 alized 0.15 Normalized 1637 0.10 0.10 1.5 0.05 1638 5 2 1.0 ອັ້ 1.0 0.05 ş 0.00 1639 0.5 0.5 0.00 -0.05 -0.05 0.0 -0.10 0.0 1640 100 200 300 400 500 ò 100 200 300 400 500 100 200 300 400 500 ò 100 200 300 400 Gradient steps (×10³) Gradient steps (×103) Gradient steps (×103) Gradient steps (×103) 1641 Safety-Gymnasium-PointButton1 Safety-Gymnasium-PointButton2 0.15 4.0 0.15 4.0 1642 3.5 3.5 0.5 st 1643 0.10 0.10 3.0 ي 2.5 2.5 1644 Re 0.05 Re 0.05 2.0 g 2.0 alized 1645 0.00 1.5 0.00 1.5 1646 ອັ້ 1.0 ອັ້ 1.0 Por ş -0.05 .n n 05 05 1647 0.0 0.0 -0.10-0.10 1648 ò 100 200 300 400 500 ò 100 200 300 400 500 ò 100 200 300 400 500 Ó 100 200 300 400 500 Gradient steps (×103) Gradient steps (×10³) Gradient steps (×103) Gradient steps (×10³ 1649

between SDQC and FISOR, we plot the training curves of both algorithms on the DSRL benchmark
(Liu et al., 2023a) in Figure 11-14. The experimental results indicate that SDQC exhibits a higher
level of safety assurance during training and achieves higher rewards in the majority of tasks.

Figure 12: Training curves of SDQC and FISOR on the "Point" agent with tasks "Goal," "Push," and "Button" in the Safety-Gymnasium domain.



Figure 13: Training curves of SDQC and FISOR on the agents "HalfCheetah", "Ant", and "Swimmer" with velocity constraints task in the Safety-Gymnasium domain.



Figure 14: Training curves of SDQC and FISOR on the agents "Ball", "Car", "Ant", and "Drone" with tasks "Circle" and "Run" in the Bullet-Safety domain.

G LIMITATIONS AND FUTURE WORKS

One limitation of our current study arises from the substantial computational demands associated with training the SDQC model. This is particularly notable due to the necessity of executing three distinct training phases and the utilization of complex network architectures in certain scenarios. Despite this challenge, the remarkable cost-effectiveness and robustness to seed variance exhibited by our model mitigate these weaknesses. Looking ahead, our future research endeavors will prioritize the optimization of the training pipeline and the simplification of network structures to enhance training efficiency while maintaining performance standards.