

---

# ExID: Offline RL with Intuitive Expert Insights in Limited-Data Settings

---

Anonymous Author(s)

Affiliation

Address

email

## Abstract

1 With the ability to learn from static datasets, Offline Reinforcement Learning (RL)  
2 emerges as a compelling avenue for real-world applications. However, state-of-the-  
3 art offline RL algorithms perform sub-optimally when confronted with limited data  
4 confined to specific regions within the state space. The performance degradation  
5 is attributed to the inability of offline RL algorithms to learn appropriate actions  
6 for rare or unseen observations. This paper proposes a novel domain knowledge-  
7 based regularization technique and adaptively refines the initial domain knowledge  
8 to considerably boost performance in limited data with partially omitted states.  
9 The key insight is that the regularization term mitigates erroneous actions for  
10 sparse samples and unobserved states covered by domain knowledge. Empirical  
11 evaluations on standard discrete environment datasets demonstrate a substantial  
12 average performance increase compared to ensemble of domain knowledge and  
13 existing offline RL algorithms operating on limited data.

## 14 1 Introduction

15 Offline RL [9, 1], also referred to as batch RL, is a learning approach that focuses on extracting  
16 knowledge solely from static datasets. This class of algorithms has a wider range of applications being  
17 particularly appealing to real-world data sets from business [46], healthcare [25], and robotics [35].  
18 However, offline RL poses unique challenges, including over-fitting and the need for generalization  
19 to data not present in the dataset. To surpass the behavior policy, offline RL algorithms need to  
20 query Q values of actions not in the dataset, causing extrapolation errors [21]. Most offline RL  
21 algorithms address this problem by enforcing constraints that ensure that the learned policy does not  
22 deviate too far away from the data set’s state action distribution [13, 11] or is conservative towards  
23 Out-of-Distribution (OOD) actions [21, 20]. However, such approaches are designed on coherent  
24 batches [13], which do not account for OOD states.

25 In many domains, such as business and healthcare, available data is scarce and often confined to expert  
26 behaviors within a limited state space. *For example, a sales recommendation system, where historic*  
27 *data may not contain details about many active users and operator gives coupon of higher value to*  
28 *attract sales.* Learning on such limited data sets can curtail the generalization capabilities of state-of-  
29 the-art (SOTA) offline RL algorithms, resulting in sub-optimal performance [23]. We illustrate this  
30 limitation via Fig 1. In Fig 1a) the state action space of a simple Mountain Car environment [27] is  
31 plotted for an expert dataset [32] and a partial dataset with first 10% samples from the entire dataset.  
32 Fig 1b) shows the average reward obtained over these data sets and the average difference between  
33 the Q value of action taken by the under-performing Conservative Q Learning (CQL) [21] agent and  
34 the action in the full expert dataset for unseen states. It can be observed that the performance of the  
35 offline RL agent considerably drops. This is attributed to the critic overestimating the Q value of  
36 non-optimal actions for states that do not occur in the dataset while training.

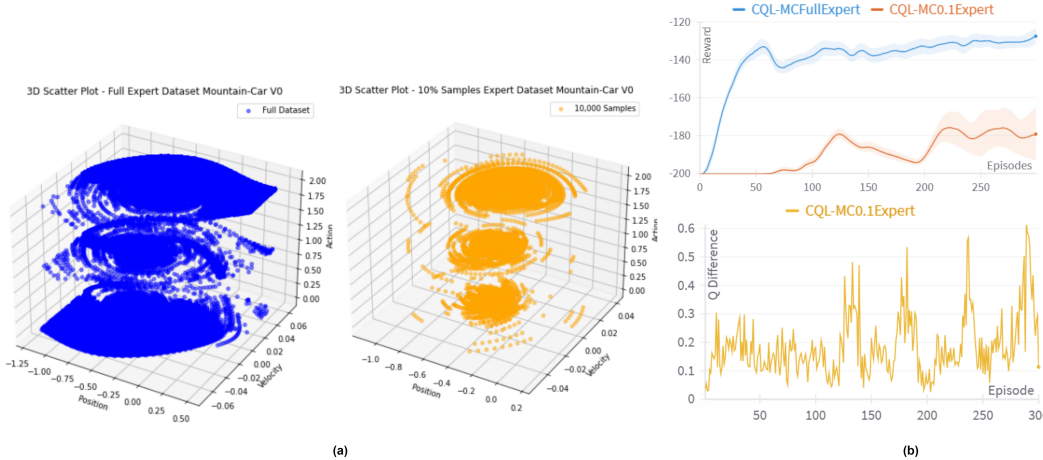


Figure 1: a) Full expert, Mountain Car dataset, and reduced dataset with first 10% samples showing distribution of state (position, velocity) and action b) CQL agent converging to a sub-optimal policy for reduced dataset exhibiting high Q values for actions different from actions in the expert dataset for unseen states.

37 In numerous real-world applications, expert insights regarding the general behavior of a policy are  
 38 often accessible [33]. For example, sales operators often distribute lower discount coupons to active  
 39 users to maximize profit. While these insights may not be optimal, they serve as valuable guidelines  
 40 for understanding the overall behavior of the policy. A rich literature in knowledge distillation [18]  
 41 has shown that teacher networks trained on domain knowledge can transfer knowledge to another  
 42 network unaware of it. This work aims to leverage a teacher network mimicking simple decision  
 43 tree-based domain knowledge to help offline RL generalize in limited data settings.

44 The paper makes the following novel contributions:

- 45 • We introduce an algorithm dubbed **ExID**, leveraging intuitive human obtainable expert  
 46 insights. The domain expertise is incorporated into a teacher policy, which improves offline  
 47 RL in limited-data settings through regularization.
- 48 • The teacher based on expected performance improvement of the offline policy during  
 49 training, improving the teacher network beyond initial heuristics.
- 50 • We demonstrate the effectiveness of our methodology on *real sales promotion dataset*,  
 51 several discrete OpenAI gym and Minigrid environments with standard offline RL data sets  
 52 and show that ExID significantly exceeds the performance when faced with limited data.

## 53 2 Related Work

54 This work improves offline RL learning on batches sampled from static datasets using domain  
 55 expertise. One of the major concerns in offline RL is the erroneous extrapolation of OOD actions  
 56 [13]. Two techniques have been studied in the literature to prevent such errors. 1) Constraining the  
 57 policy to be close to the behavior policy 2) Penalizing overly optimistic Q values [24]. We discuss a  
 58 few relevant algorithms following these principles. In Batch-Constrained deep Q-learning (BCQ)  
 59 [13] candidate actions sampled from an adversarial generative model are considered, aiming to  
 60 balance proximity to the batch while enhancing action diversity. Algorithms like Random Ensemble  
 61 Mixture Model (REM) [2], Ensemble-Diversified Actor-Critic (EDAC) [3] and Uncertainty Weighted  
 62 Actor-Critic (UWAC) [42] penalize the Q value according to uncertainty by either using Q ensemble  
 63 networks or directly weighting the loss with uncertainty. CQL [21] enforces regularization on Q-  
 64 functions by incorporating a term that reduces Q-values for OOD actions while increasing Q-values  
 65 for actions within the expected distribution. However, these algorithms do not handle OOD actions  
 66 for states not in the static knowledge and can have errors induced by changes in transition probability.  
 67 Integration of domain knowledge in offline RL, though an important avenue, has not yet been

68 extensively explored. Domain knowledge incorporation has improved online RL with tight regret  
69 bounds [33, 4]. In offline RL, bootstrapping via blending heuristics computed using Monte-Carlo  
70 returns with rewards has shown to outperform SOTA algorithms by 9% [15]. Recent works improve  
71 offline RL by incorporating a safety expert [40] and preference query [44], contrary to our work  
72 which improves imperfect domain knowledge. The closest to our work is Domain Knowledge guided  
73 Q learning (DKQ) [46] where domain knowledge is represented in terms of action importance and  
74 the Q value is weighted according to importance. However, obtaining action importance in practical  
75 scenarios is nontrivial.

### 76 3 Preliminaries

77 A DRL setting is represented by a Markov Decision Process (MDP) formalized as  $(S, A, T, r, \rho_0, \gamma)$ .  
78 Here,  $S$  denotes the state space,  $A$  signifies the action space,  $T(s'|s, a)$  represents the transition prob-  
79 ability distribution,  $r : S \times A \rightarrow \mathbb{R}$  is the reward function,  $\rho_0$  represents the initial state distribution,  
80 and  $\gamma \in (0, 1]$  is the discount factor. The primary objective of any DRL algorithm is to identify an  
81 optimal policy  $\pi(a|s)$  that maximizes  $\mathbb{E}_{s_t, a_t} [\sum_{t=0}^{\infty} \gamma^t r(s_t, a_t)]$  where,  $s_0 \sim d_0(\cdot)$ ,  $a_t \sim \pi(\cdot|s_t)$ , and  
82  $s' \sim T(\cdot|s_t, a_t)$ . Deep Q networks (DQNs) [26] learn this objective by minimizing the Bellman resid-  
83 ual  $(Q_\theta(s, a) - B^{\pi_\theta} Q_\theta(s, a))^2$  where  $B^{\pi_\theta} Q_\theta(s, a) = \mathbb{E}_{s' \sim T} [r(s, a) + \gamma \mathbb{E}_{a' \sim \pi_\theta(\cdot|s')} [Q_{\theta'}(s', a')]]$ .  
84 The policy  $\pi_\theta$  chooses actions that maximize the Q value  $\max_{a' \in A} Q_\theta(s', a')$ . However, in offline  
85 RL where transitions are sampled from a pre-collected dataset  $\mathcal{B}$ , the chosen action  $a'$  may exhibit a  
86 bias towards OOD actions with inaccurately high Q-values. To handle the erroneous propagation  
87 from OOD actions, CQL [22] learns conservative Q values by penalizing OOD actions. The CQL  
88 loss for discrete action space is given by

$$\begin{aligned} \mathcal{L}_{cql}(\theta) = \min_Q \alpha \mathbb{E}_{s \sim \mathcal{B}} [\log \sum_a \exp(Q_\theta(s, a))] - \\ \mathbb{E}_{a \sim \mathcal{B}|s} [Q_\theta(s, a)] + \frac{1}{2} \mathbb{E}_{s, a, s' \sim \mathcal{B}} [Q_\theta - Q_{\theta'}]^2 \end{aligned} \quad (1)$$

89 Eq. 1 encourages the policy to be close to the actions seen in the dataset. However, CQL works on the  
90 assumption of coherent batches, i.e., if  $(s, a, s') \in \mathcal{B}$ , then  $s' \in \mathcal{B}$ . There is no provision for handling  
91 OOD actions for  $s \notin \mathcal{B}$ , which can lead to policy failure when data is limited. In the next sections, we  
92 present ExID, a domain knowledge-based approach to improve performance in data-scarce scenarios.

### 93 4 Problem Setting and Methodology

94 In our problem setting, the RL agent learns the policy on a limited dataset with rare and unseen  
95 demonstrations. We define the characteristics of this dataset as follows:

96 **Definition 4.1.** A reduced buffer  $\mathcal{B}_r$  is a proper subset of the full dataset  $\mathcal{B}$  i.e.,  $\mathcal{B}_r \subset \mathcal{B}$  satisfying  
97 the following conditions:

- 98 • Some states in  $\mathcal{B}$  are not present in  $\mathcal{B}_r$ , i.e.,  $\exists s' \in \mathcal{B} \wedge \forall (s, a, s') : (s, a, s') \notin \mathcal{B}_r$
- 99 • The number of samples  $N(s, a, s')$  for some transitions in  $\mathcal{B}$  are less in  $\mathcal{B}_r$  i.e.,  $\exists (s, a, s') \in$   
100  $\mathcal{B} : N(s, a, s')_{\mathcal{B}_r} < N(s, a, s')_{\mathcal{B}}$

101 We observe, performing *Q-Learning by sampling from a limited buffer  $\mathcal{B}_r$  may not converge*  
102 *to an optimal policy for the MDP  $M_{\mathcal{B}}$  representing the full buffer.* This can be shown as a special  
103 case of (Theorem 1, [13]) as  $p_{\mathcal{B}}(s'|s, a) \neq p_{\mathcal{B}_r}(s'|s, a)$  and no Q updates for  $(s, a) \notin \mathcal{B}_r$  leading to  
104 sub-optimal policy. Please refer to the App. B for analysis and example.

105 We also assume a set of common sense rules in the form of domain knowledge is available. Domain  
106 knowledge  $\mathcal{D}$  is defined as hierarchical decision nodes capturing  $S \rightarrow A$  as represented by Eq. 2.  
107 Each decision node  $T_{\eta_i}$  is represented by a constraint  $\phi_{\eta_i}$  and Boolean indicator  $\mu_{\eta_i}$  function selects  
108 the branch to be traversed based on  $\phi_{\eta_i}$ .

$$\begin{aligned} \text{Action} &= \begin{cases} a_{\eta_i} & \text{if leaf} \\ \mu_{\eta_i} T_{\eta_i \swarrow}(s) + (1 - \mu_{\eta_i}) T_{\eta_i \searrow}(s) & \text{o/w} \end{cases} \\ \mu_{\eta_i}(s) &= \begin{cases} 1 & \text{if } s \models \phi_{\eta_i} \\ 0 & \text{o/w} \end{cases} \end{aligned} \quad (2)$$

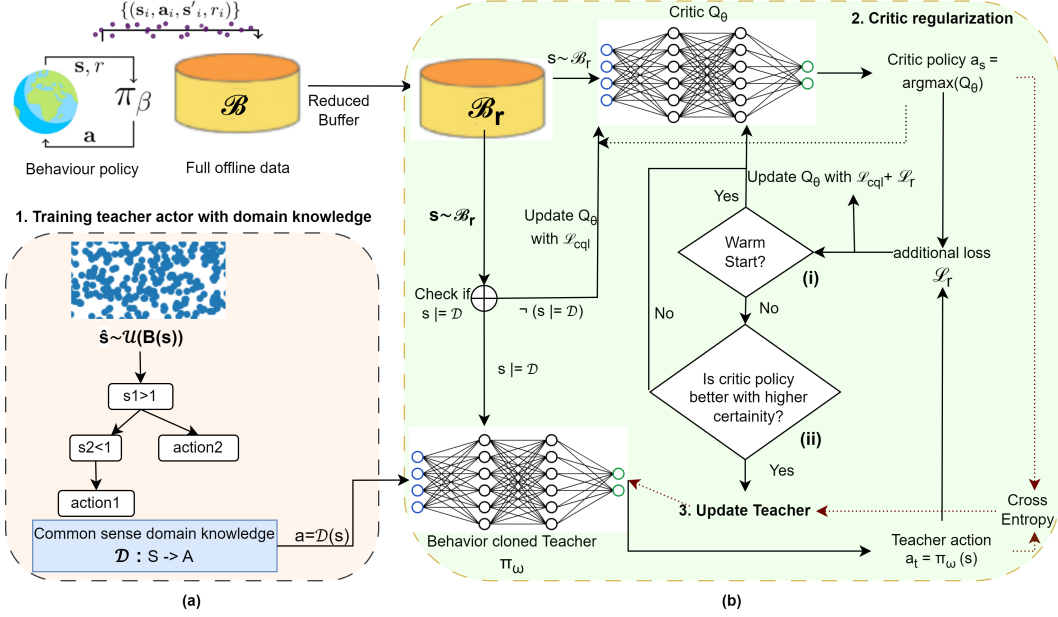


Figure 2: Overview of the proposed methodology (a) Training a teacher policy network with domain knowledge and synthetic data (b) Updating the offline RL critic network with teacher network

109 We assume that  $\mathcal{D}$  gives heuristically reasonable actions for  $s \models \mathcal{D}$  and  $S_{\mathcal{D}} \cap S_{\mathcal{B}_r} \neq \emptyset$  where  $S_{\mathcal{D}}, S_{\mathcal{B}_r}$   
 110 are the state coverage of  $\mathcal{D}$  and  $\mathcal{B}_r$ .

111 **Training Teacher:** An overview of our methodology is depicted in Fig 2. We first construct a  
 112 trainable actor network  $\pi_t^\omega$  parameterized by  $\omega$  from  $\mathcal{D}$ , Fig 2 step 1. For training  $\pi_t^\omega$  synthetic  
 113 data  $\hat{S}$  is generated by sampling states from a uniform random distribution over state boundaries  
 114  $B(s)$ ,  $\hat{S} = \mathcal{U}(B(S))$ . Note that this does not represent the true state distribution and may have state  
 115 combinations that will never occur. We train  $\pi_t^\omega$  using behavior cloning where state  $\hat{s} \sim \hat{S}$  is checked  
 116 with root decision node in Eq. 2. A random action is chosen if  $\hat{s}$  does not satisfy decision node  $T_{\eta_0}$   
 117 or leaf action is absent. If  $\hat{s}$  satisfies a  $T_{\eta_i}$ ,  $T_{\eta_i}$  is traversed and action  $a_{\eta_i}$  is returned from the  
 118 leaf node. This is illustrated in Fig 2 (a). We term the pre-trained actor network  $\pi_t^\omega$  as the teacher policy.

119 **Regularizing Critic:** We now introduce Algo 1 (App C) to train an offline RL agent on  $\mathcal{B}_r$ . Algo 1  
 120 takes  $\mathcal{B}_r$  and pretrained  $\pi_t^\omega$  as input. The algorithm uses two hyper-parameters, warm start parameter  
 121  $k$  and mixing parameter  $\lambda$ . A critic network  $Q_s^\theta$  with Monte-Carlo (MC) dropout and target network  
 122  $Q_s^{\theta'}$  are initialized. ExID is divided into two phases. In the first phase, we aim to warm start the critic  
 123 network  $Q_s^\theta$  with actions from  $\pi_t^\omega$  as shown in Fig 2b (i). However, this must be done selectively  
 124 as the teacher's policy is random around the states that do not satisfy domain knowledge. In each  
 125 iteration, we first check the states sampled from a mini-batch of  $\mathcal{B}_r$  with  $\mathcal{D}$ . For the states which  
 126 satisfy  $\mathcal{D}$  we compute the teacher action  $\pi_t^\omega(s)$  and critic's action  $\text{argmax}_a(Q_s^\theta(s, a))$  and collect it  
 127 in lists  $a_t, a_s$ , Algo 1 lines 4-10. Our main objective is to keep actions chosen by the critic network  
 128 for  $s \models \mathcal{D}$  close to the teacher's policy. To achieve this, we introduce a regularization term:

$$\mathcal{L}_r(\theta) = \underbrace{\mathbb{E}_{s \sim \mathcal{B}_r \wedge s \models \mathcal{D}}}_{\text{states matching domain rule}} \underbrace{[Q_s^\theta(s, a_s) - Q_s^\theta(s, a_t)]^2}_{\text{Q regularizer}} \quad (3)$$

129 Eq 3 incentivizes the critic to increase Q values for actions from  $\pi_t^\omega$  and decreases Q values for other  
 130 actions when  $\text{argmax}_a(Q_s^\theta(s, a)) \neq \pi_t^\omega(s)$  for states that satisfy domain knowledge. Note that Eq 3  
 131 will only be 0 when  $\text{argmax}_a(Q_s^\theta(s, a)) = \pi_t^\omega(s)$  for  $s \models \mathcal{D}$ . It is also set to 0 for  $s \not\models \mathcal{D}$ . However,  
 132 since  $\pi_t^\omega$  mimicking heuristic rules is sub-optimal, it is also important to incorporate learning from  
 133 the data. The final loss is a combination of Eq. 1 and Eq. 3 with a mixing parameter  $\lambda \in [0, 1]$ :

$$\mathcal{L}(\theta) = \mathcal{L}_{cqi}(\theta) + \lambda \mathbb{E}_{s \sim \mathcal{B}_r \wedge s \models \mathcal{D}} [Q_s^\theta(s, a_s) - Q_s^\theta(s, a_t)]^2 \quad (4)$$

134 The choice of  $\lambda$  and the warm start parameter  $k$  depends on the quality of  $\mathcal{D}$ . In the case of perfect  
 135 domain knowledge,  $\lambda$  would be set to 1, and setting  $\lambda$  to 0 would lead to the vanilla CQL loss. Mixing  
 136 both the losses allows the critic to learn both from the data in  $\mathcal{B}_r$  and knowledge in  $\mathcal{D}$ .

137 **Updating Teacher:** Given a reasonable warm start, the critic is expected to give higher Q values  
 138 for optimal actions for  $s \in \mathcal{D} \cap \mathcal{B}_r$  as it learns from data. We aim to leverage this knowledge  
 139 to enhance the initial teacher policy  $\pi_t^\omega$  trained on heuristic domain knowledge. For  $s \sim \mathcal{B}$  and  
 140  $s \models \mathcal{D}$ , we calculate the average Q values over critic actions and teacher actions and check which  
 141 one is higher in Algo 1 line 11 which refers to Cond. 6. For brevity  $\mathbb{E}_{s \sim \mathcal{B}_r \wedge s \models \mathcal{D}}$  is written as  $\mathbb{E}$ .  
 142 If  $\mathbb{E}(Q_s^\theta(s, a_s)) > \mathbb{E}(Q_s^\theta(s, a_t))$  it denotes the critic expects a better return on an average over its  
 143 own policy than the teacher’s policy. Hence, we can use the critic’s policy to update  $\pi_t^\omega$ , making  
 144 it better over  $\mathcal{D}$ . However, only checking the critic’s value can be erroneous as the critic can have  
 145 high values for OOD actions. We check the average uncertainty of the predicted Q values to prevent  
 146 the teacher from getting updated by OOD actions. Uncertainty has been shown to be a good metric  
 147 for OOD action detection by [42, 3]. A well-established methodology to capture uncertainty is  
 148 predictive variance, which takes inspiration from Bayesian formulation for the critic function and  
 149 aims to maximize  $p(\theta|X, Y) = p(Y|X, \theta)p(\theta)/p(Y|X)$  where  $X = (s, a)$  and  $Y$  represents the true  
 150 Q value of the states. However,  $p(Y|X)$  is generally intractable and is approximated using Monte  
 151 Carlo (MC) dropout, which involves including dropout before every layer of the critic network and  
 152 using it during inference [14]. Following [42], we measure the uncertainty of prediction using Eq 5.

$$Var^T [Q(s, a)] \approx \frac{1}{T} \sum_{t=1}^T [Q(s, a) - \bar{Q}(s, a)]^2 \quad (5)$$

153 Eq 5 estimates the variance of Q value  $Q(s, a)$  for an action  $a$  using  $T$  forward passes on the  $Q_s^\theta(s, a)$   
 154 with dropout where  $\bar{Q}(s, a)$  represents the predictive mean. We check the average uncertainty of  
 155 the Q value for action chosen by the critic and teacher policy over the states that match domain  
 156 knowledge in a batch. The teacher network is updated using the critic’s action only when the policy  
 157 expects a higher average Q return on its action and the average uncertainty of taking this action is  
 158 lower than the teacher action.  $\mathbb{E}(Var^T Q_s^\theta(s_r, a_s)) < \mathbb{E}(Var^T Q_s^\theta(s_r, a_t))$  indicates the actions were  
 159 learned from the expert data in the buffer and are not OOD samples. The condition is summarized in  
 160 cond. 6:

$$\mathbb{E}(Q_s^\theta(s_r, a_s)) > \mathbb{E}(Q_s^\theta(s_r, a_t)) \wedge \mathbb{E}(Var^T Q_s^\theta(s_r, a_s)) < \mathbb{E}(Var^T Q_s^\theta(s_r, a_t)) \quad (6)$$

161 We update the teacher with cross-entropy described in Eq 7:

$$\mathcal{L}(\omega) = - \sum_{s \models \mathcal{D}} (\pi_t^\omega(s) \log(\pi_s(s))) \quad (7)$$

162 where,  $\pi_s(s, a) = \frac{e^{Q(s, a)}}{\sum_{a'} Q(s, a')}$ . When the critic’s policy is better than the teacher’s policy,  $\mathcal{L}_r(\theta)$  is  
 163 set to 0 Algo 1 Lines 11 to 13. Finally, the critic network is updated using calculated loss  $\mathcal{L}(\theta)$  Algo  
 164 1 Lines 17-18. We theoretically analyse the implications of using ExID in propositions 4.2 and 4.3.

165 **Proposition 4.2.** Denote  $\hat{\pi}$  as the policy learned by ExID,  $\pi_u$  as any offline RL policy learned on  $\mathcal{B}_r$   
 166 and optimal Q function as  $Q^*$  and V function as  $V^*$ . Then it holds that

$$\eta(\hat{\pi}) - \eta(\pi_u) \geq \mathbb{E}_{s \sim O | \pi_u} [V^*(s) - Q^*(s, \pi_u(s))] - \bar{\rho}_{\hat{\pi}} \alpha$$

167 Where  $\alpha = \mathbb{E}_{s \sim O} [V^*(s) - Q^*(s, \hat{\pi}(s))]$ ,  $\bar{\rho}_{\hat{\pi}}(s) = [\frac{1}{|S_{\hat{\pi}}|(1-\gamma)}, \frac{1}{1-\gamma}]$  ( $|S_{\hat{\pi}}|$  is the number of different  
 168 states observed by  $\hat{\pi}$ ) and  $O \notin \mathcal{B}_r$ . Here  $\alpha$  denotes the quality of regularized action for  $s \notin \mathcal{B}_r$ . Hence,  
 169 updating  $\pi_t^\omega$  is important as high divergence of action from the optimal can lead to performance  
 170 degradation. In offline RL, the extrapolation error for non optimal action is usually high for states not  
 171 observed in dataset (as illustrated in 1b), regularization can lead performance improvement when  $\pi_t^\omega$   
 172 is reasonable. Furthermore, in ExID coarse actions from  $\pi_t^\omega$  are updated driving them closer to the  
 173 optimal actions, improving the performance lower bound. Additionally  $\pi_t^\omega$  increases  $|S_{\hat{\pi}}|$  making  
 174  $\bar{\rho}_{\hat{\pi}} \ll 1$  in practice further improving the performance lower bound. *Proof is deferred to App. A.*

175 **Proposition 4.3.** *ExID reduces generalization error if  $Q^*(s, \pi_t^\omega(s)) > Q^*(s, \pi_u(s))$  for  $s \in \mathcal{D} \cap \mathcal{B}_r$ .*  
176 *Proof is deferred to App. A. In the next section, we discuss our empirical evaluations.*

## 177 5 Empirical Evaluations

178 We investigate the following through our empirical evaluations: 1. *Does ExID perform better than*  
179 *combining  $\mathcal{D}$  and offline RL algos on different environments with datasets exhibiting rare and OOD*  
180 *states Sec 5.2?* 2. *Does ExID generalize to OOD states covered by  $\mathcal{D}$  Sec 5.4?* 3. *What is the effect of*  
181 *varying  $k$ ,  $\lambda$  and updating  $\pi_t^\omega$  Sec 5.5?* 4. *How does performance vary with the quality of  $\mathcal{D}$  Sec 5.6?*

### 182 5.1 Experimental Setting

183 We evaluate our methodology on open-AI gym [5], MiniGrid [6] and *real sales promotion (SP)* [30]  
184 offline data sets. All our data sets are generated using standard methodologies defined in [32, 31]  
185 *except SP which is generated by human operators.* All experiments have been conducted on a  
186 Ubuntu 22.04.2 LTS system with 1 NVIDIA K80 GPU, 4 CPUs, and 61GiB RAM. App. F notes the  
187 hyperparameter values and network architectures.

188 **Dataset:** We experiment on three types of data sets. *Expert Data-set* [10, 16, 22] generated using  
189 an optimal policy without any exploration with high trajectory quality but low state action coverage.  
190 *Replay Data-set* [2, 13] generated from a policy while training it online, exhibiting a mixture of  
191 multiple behavioral policies with high trajectory quality and state action coverage. *Noisy Data-set*  
192 [12, 13, 22, 16] generated using an optimal policy that also selects random actions with  $\epsilon$  greedy  
193 strategy where  $\epsilon = 0.2$  having low trajectory quality and high state action coverage. Additionally we  
194 also experiment on human generated dataset for sales promotion task.

195 **Baselines:** We do comparative studies on 10 baselines for OpenAI gym datasets. The first baseline  
196 simply checks the conditions of  $\mathcal{D}$  and applies corresponding actions in execution. The performance  
197 of this baseline shows that  $\mathcal{D}$  is imperfect and does not achieve the optimal reward. CQL SE is  
198 from [40] where the expert is replaced by  $\mathcal{D}$ . The other baselines are an ensemble of  $\mathcal{D}$  and eight  
199 algorithms popular in the Offline RL literature for discrete environments. These algorithms include  
200 Behavior Cloning (BC) [29], Behaviour Value Estimation (BVE) [16], Quantile Regression DQN  
201 (QRDQN) [7], REM, MCE, BCQ, CQL and Critic Regularized Regression Q-Learning (CRR) [41].  
202 *For a fair comparison, we use actions from domain knowledge for states not in the buffer and actions*  
203 *from the trained policy for other states to obtain the final reward.* Hence, each algorithm is renamed  
204 with the suffix D in Table 5.1.

205 **Limiting Data:** To create limited-data settings for benchmark datasets, we first extract a small  
206 percentage of samples from the full dataset and remove some of the samples based on state conditions.  
207 This is done to ensure the reduced buffer satisfies the conditions defined in Def 4.1. We describe  
208 the specific conditions of removal in the next section. Further insights and the state visualizations  
209 for selected reduced datasets are in App H. **Note : no data reduction has been performed on SP**  
210 **dataset to demonstrate a real dataset exhibits characteristics of reduced buffer.**

### 211 5.2 Performance across Different Datasets

212 Our results for OpenAI gym environments are summarised in Table 5.1 and Minigrid in Table 3 (App  
213 D). We observe the performance of offline RL algorithms degrades substantially when part of the data  
214 is not seen and trajectory ratios change. For these cases with only 10% partial data, ExID surpasses  
215 the performance by at least 27% in the presence of reasonable domain knowledge. The proposed  
216 method performs strongest on the replay dataset where the contribution of  $L_r(\theta)$  is significant due  
217 to state coverage, and the teacher learns from high-quality trajectories. Environment details are  
218 described in the App. D. All domain knowledge trees are shown in the App. D Fig 10. We describe  
219 limiting data conditions and domain knowledge specific to the environment as follows:

220 **Mountain Car Environment:** [27] We use simple, intuitive domain knowledge in this environment  
221 shown in the App. D Fig 10 (c), which represents taking a left action when the car is at the bottom of  
222 the valley with low velocity to gain momentum; otherwise, taking the right action to drive the car up.  
223 Fig 6 (c) shows the state action pairs this rule generates on states sampled from a uniform random  
224 distribution over the state boundaries. It can be observed that the states of  $\mathcal{D}$  cover part of the missing

Table 1: Average reward  $\uparrow$  obtained during online evaluation over 3 seeds on openAI gym envs

ENV DATA	DATA TYPE	$\mathcal{D}$	QRDQN D	REM D	BVE D	CRR D	MCE D	BC D	BCQ D	CQL D	CQL SE	CQL (FULL)	EXID (OURS)
MOUNTAIN CAR	EXPERT		-168.2	-147.7	-175.36	-157.2	-152	-181.38	-172.9	-167.49	-161.33	-128.63	-125.5
		$\pm$	$\pm$	$\pm$	$\pm$	$\pm$	$\pm$	$\pm$	$\pm$	$\pm$	$\pm$	$\pm$	$\pm$
	REPLAY	-159.9	33.71	21.54	25.16	39.09	37.41	28.60	27.5	12.3	18.57	10.94	2.60
		$\pm$	$\pm$	$\pm$	$\pm$	$\pm$	$\pm$	$\pm$	$\pm$	$\pm$	$\pm$	$\pm$	$\pm$
	NOISY	52.28	-137.14	-136.26	-152.0	-137.23	-139.91	-137.26	-136.29	-140.38	-150.67	-135.4	-105.79
		$\pm$	$\pm$	$\pm$	$\pm$	$\pm$	$\pm$	$\pm$	$\pm$	$\pm$	$\pm$	$\pm$	$\pm$
CART POLE	EXPERT		39.27	40.15	35.06	42.79	40.01	43.04	36.15	33.58	16.68	3.74	11.38
		$\pm$	$\pm$	$\pm$	$\pm$	$\pm$	$\pm$	$\pm$	$\pm$	$\pm$	$\pm$	$\pm$	$\pm$
	REPLAY	57.0	-141.61	-134.99	-173.95	-178.99	-168.69	-140.0	-144.52	-179.8	-126.96	-107.06	-109.9
		$\pm$	$\pm$	$\pm$	$\pm$	$\pm$	$\pm$	$\pm$	$\pm$	$\pm$	$\pm$	$\pm$	$\pm$
	NOISY	5.35	33.04	32.60	39.60	23.58	38.78	28.5	43.04	29.99	17.84	12.73	13.45
		$\pm$	$\pm$	$\pm$	$\pm$	$\pm$	$\pm$	$\pm$	$\pm$	$\pm$	$\pm$	$\pm$	$\pm$
LUNAR LANDER	EXPERT		33.23	41.31	16.16	15.24	16.1	225.76	165.36	121.8	155.78	364.1	307.18
		$\pm$	$\pm$	$\pm$	$\pm$	$\pm$	$\pm$	$\pm$	$\pm$	$\pm$	$\pm$	$\pm$	$\pm$
	REPLAY	52.48	3.17	8.76	9.41	5.62	4.4	74.39	15.01	14.0	26.47	22.15	137.72
		$\pm$	$\pm$	$\pm$	$\pm$	$\pm$	$\pm$	$\pm$	$\pm$	$\pm$	$\pm$	$\pm$	$\pm$
	NOISY	26.51	149.09	180.70	11.1	11.24	9.16	144.43	144.76	131.97	113.37	250.02	340.26
		$\pm$	$\pm$	$\pm$	$\pm$	$\pm$	$\pm$	$\pm$	$\pm$	$\pm$	$\pm$	$\pm$	$\pm$
			14.05	62.79	2.13	2.71	0.25	2.41	6.04	23.23	5.88	55.02	30.58
			$\pm$	$\pm$	$\pm$	$\pm$	$\pm$	$\pm$	$\pm$	$\pm$	$\pm$	$\pm$	$\pm$
			6.40	0.58	3.77	7.49	2.04	14.67	14.08	22.05	22.05	37.79	38.64
			$\pm$	$\pm$	$\pm$	$\pm$	$\pm$	$\pm$	$\pm$	$\pm$	$\pm$	$\pm$	$\pm$
			5.14	-184.84	-681.67	8.79	19.71	38.40	-45.99	65.43	53.22	167.74	161.34
			$\pm$	$\pm$	$\pm$	$\pm$	$\pm$	$\pm$	$\pm$	$\pm$	$\pm$	$\pm$	$\pm$
			25.10	26.45	34.86	25.38	10.52	23.21	30.47	71.37	78.85	29.4	17.10
			$\pm$	$\pm$	$\pm$	$\pm$	$\pm$	$\pm$	$\pm$	$\pm$	$\pm$	$\pm$	$\pm$
			-444.20	-556.81	-572	-131.21	-115.23	136.63	111.47	61.83	87.70	187.72	156.03
			$\pm$	$\pm$	$\pm$	$\pm$	$\pm$	$\pm$	$\pm$	$\pm$	$\pm$	$\pm$	$\pm$
			12.20	21.39	27.93	31.97	18.16	12.40	54.67	45.57	18.20	25.62	56.67
			$\pm$	$\pm$	$\pm$	$\pm$	$\pm$	$\pm$	$\pm$	$\pm$	$\pm$	$\pm$	$\pm$
			-4.81	21.41	28.65	-158.27	-50.47	98.62	101.59	5.01	40.35	111	163.57
			$\pm$	$\pm$	$\pm$	$\pm$	$\pm$	$\pm$	$\pm$	$\pm$	$\pm$	$\pm$	$\pm$
			97.28	14.71	12.26	7.71	15.78	28.01	30.83	128.63	65.72	52.32	49.24
			$\pm$	$\pm$	$\pm$	$\pm$	$\pm$	$\pm$	$\pm$	$\pm$	$\pm$	$\pm$	$\pm$

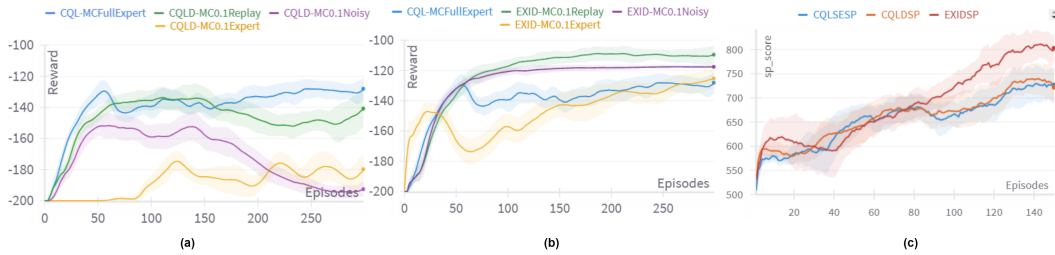


Figure 3: Performance of (a) CQL and (b) EXID on all datasets for Mountain Car during online evaluation (c) Evaluation curves for the sales promotion dataset

225 data in Fig 1 (a). For limiting datasets, we remove states with position  $> -0.8$ . The performance of  
 226 CQLD and ExID are shown in Fig 3 (a),(b) where ExID surpasses CQLD for all three datasets.

227 **Cart-pole Environment:** For this environment, we use domain knowledge from [33], which aims to  
 228 move in the direction opposite to the lean of the pole, keeping the cart close enough to the center. If  
 229 the cart is close to an edge, the domain knowledge attempts to account for the cart’s velocity and  
 230 recenter the cart. The full tree is given in the App. D Fig 10 (a). We remove states with cart velocity  
 231  $> -1.5$  to create the reduced buffer.

232 **Lunar-Lander Environment:** We borrow the decision nodes from [34] and get actions from a  
 233 sub-optimal policy trained online with an average reward of 52.48. The full set of decision nodes is  
 234 shown in the App. D Fig 10 (b).  $\mathcal{D}$  focuses on keeping the lander balanced when the lander is above  
 235 ground. When the lander is near the surface,  $\mathcal{D}$  focuses on keeping the y velocity lower. To create the  
 236 reduced datasets, we remove data of lander angle  $< -0.04$ .

237 **Mini-Grid Environments:** For our experiments, we choose two environments: Random Dynamic  
 238 Obstacles 6X6 and LavaGapS 7X7. We use intuitive domain knowledge which avoids crashing into  
 239 obstacles in front, left, or right of agent ref. App. D Fig 10 (d), (e). We remove states with obstacles  
 240 on the right for creating limited data settings. Due to limitation of space we report the results of the  
 241 best-performing algorithms on the replay dataset in Table 3 (App D).

### 242 5.3 Case study on real human generated Sales Promotion (SP) dataset

243 SP dataset and environment [30] simulates a real-world sales promotion platform. The number of  
 244 coupons and the discount the user received will affect his behavior. A higher discount will promote

245 the sales, but the cost will also increase. The goal for the platform operator is to maximize the  
 246 total profit. The horizon of the dataset is 50 days for the training and 30 days for the test. Domain  
 247 knowledge ([30], App A) : Active users can be given more coupons with lower discount to maximize  
 248 profit. We model this as  $order\_number > 60 \wedge Avg_{fee} > 0.8 \implies [5, 0.95]$  where action 1 is number  
 249 of coupons range  $[0, 5]$  and action 2 is coupon value (discount value =  $(1 - coupon\ value)$ ) range  
 250  $[0.6, 0.95]$ . The dataset exhibits the properties in Def 4.1 as first 50 days of sales does not contain  
 251 many active users as reported in the coverage column of Tab 2 depicting scarcity. The domain rule is  
 252 imperfect as coupon value and number depend on multiple factors such as user purchase history and  
 253 behavior. As illustrated in the table 2 and Fig 3 (c) the intuitive domain rule enhances performance  
 254 by 10.49% in the real dataset.

Table 2: Results on human generated Sales Promotion dataset

Dataset	$\mathcal{D}$	coverage $\mathcal{D}$	CQL + $\mathcal{D}$	CQLSE	EXID	Performance gain
Sales	654.68	20.32%	$722.06 \pm 71.40$	$727.03 \pm$	802.91	10.49%
Promotion	$\pm 20.06$			49.56	$\pm 41.69$	

#### 255 5.4 Generalization to OOD states and contribution of $\mathcal{L}_r(\theta)$

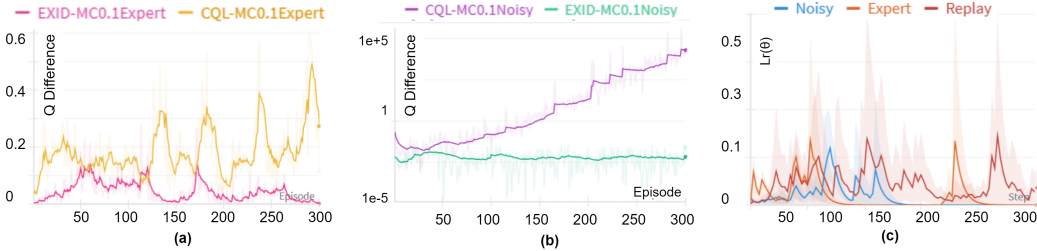


Figure 4: Q value difference between CQL and EXID for expert and policy action on states not present in the buffer for a) expert b) noisy in log scale c) contribution of  $\mathcal{L}_r(\theta)$

256 In Fig 4 (a), (b), we plot  $Q_s^\theta(s, a_{expert}) - Q_s^\theta(s, a_\theta)$  for CQL and EXID policies for different datasets  
 257 of Mountain-Car environments. Action  $a_{expert}$  is obtained from the full expert dataset where position  
 258  $> -0.8$ . We observe that the Q value for actions of CQL policy diverges from the expert policy  
 259 actions with high values for the states not in the reduced buffer, whereas ExID stays close to the  
 260 expert actions for the unseen states. This empirically shows generalization to OOD states not in the  
 261 dataset but covered by domain knowledge. In Fig 4 (d), we plot the contribution by  $\mathcal{L}_r(\theta)$  during the  
 262 training and observe the contribution is higher for replay data sets with more state coverage.

#### 263 5.5 Performance on varying $\lambda$ , $k$ , and ablation of $\pi_t^\omega$

264 We study the effect of varying  $\lambda$  on the algorithm for the given domain knowledge. We empirically  
 265 observe setting a high or a low  $\lambda$  can yield sub-optimal performance, and  $\lambda = 0.5$  generally gives  
 266 good performance. In Fig 5 (a), we show this effect for LunarLander. Plots for other environments  
 267 are in the App. G Fig 11. For  $k$  we observe setting the warm start parameter to 0 yields a sub-optimal  
 268 policy, as the critic may update  $\pi_t^\omega$  without completely learning from it. The starting performance  
 269 increases with an increase in  $k$  as shown in Fig 5 (b) for LunarLander.  $k = 30$  works best according  
 270 to empirical evaluations. Plots for other environments are in the App. G Fig 12. We show two  
 271 ablations for Cart-pole in Fig 5 (c) with no teacher update after the warm start and no inclusion of  
 272  $\mathcal{L}_r(\theta)$  after the warm start. The warm start in this environment is set to 30 episodes. Fig 5 c) shows  
 273 without teacher updated, the sub-optimal teacher drags down the performance of the policy beyond  
 274 the warm start, exhibiting the necessity of  $\pi_t^\omega$  update. Also, the student converges to a sub-optimal  
 275 policy if no  $\mathcal{L}_r(\theta)$  is included beyond the warm start.



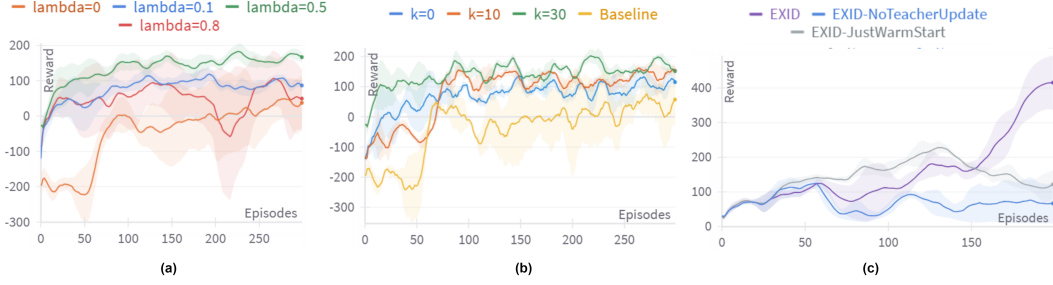


Figure 5: (a) Effect of different  $\lambda$  on the performance of ExID on Lunar Lander (b) Effect of different  $k$  on the performance of EXID on Lunar Lander (c) Performance of EXID with teacher update, no teacher update, and just warm start on Cart-pole.

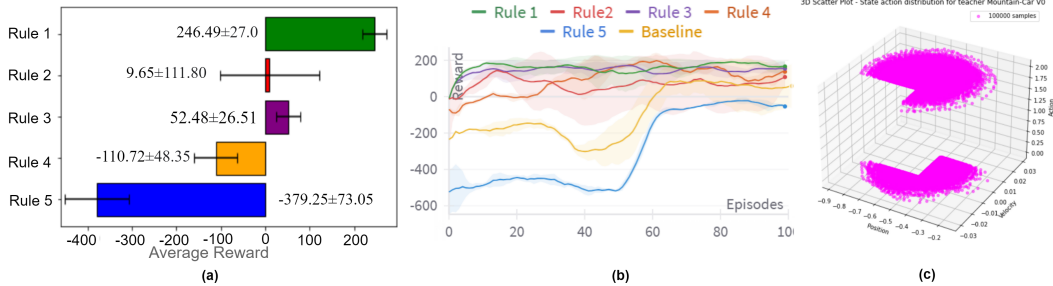


Figure 6: (a)  $\mathcal{D}$  with different average rewards (b) Performance effect on Lunar-lander (c) State distribution generated for training the teacher network for mountain-car

276 **5.6 Effect of varying  $\mathcal{D}$  quality**

277 We show the effect of choosing policies as  $\mathcal{D}$  with different average rewards for Lunar-Lander expert  
 278 data in Fig 6 (a) and (b). Rule 1 is optimal and has almost the same effect as Rule 3, which is the  $\mathcal{D}$   
 279 used in our experiments exhibiting that updating a sub-optimal  $\mathcal{D}$  can lead to equivalent performance  
 280 as optimal  $\mathcal{D}$ . Using a rule with high uncertainty, as Rule 2, induces high uncertainty in the learned  
 281 policy but performs slightly better than the baseline. Rule 4, which has a lower average reward, also  
 282 causes gains on average performance with slower convergence. Finally, Rule 5, with very bad actions,  
 283 affects policy performance adversely and leads to a performance lower than baseline CQL.

284 **6 Conclusion and Limitation**

285 In this paper, we study the effect of limited and partial data on offline RL and observe that the  
 286 performance of SOTA offline RL algorithms is sub-optimal in such settings. The paper proposes a  
 287 methodology to handle offline RL’s performance degradation using domain insights. We incorporate  
 288 a regularization loss in the CQL training using a teacher policy and refine the initial teacher policy  
 289 while training. We show that incorporating reasonable domain knowledge in offline RL enhances  
 290 performance, achieving a performance close to full data. However, this method is limited by the  
 291 quality of the domain knowledge and the overlap between domain knowledge states and reduced  
 292 buffer data. The study is also limited to discrete domains. In the future, the authors would like to  
 293 improve on capturing domain knowledge into the policy network without dependence on data and  
 294 extending the methodology to algorithms that handle continuous action space.

295 **7 Broader Impact**

296 During the trial-and-error training phase, RL agents may exhibit irrational behavior, which can be  
 297 risky and costly in real-world scenarios. As a more practical alternative to online RL, offline RL

298 utilizes pre-existing collected data to eliminate the need for real-time interactions during training.  
299 However, a drawback of offline RL is its dependence on the quality and quantity of historical data,  
300 which, when sub-optimal, could adversely affect overall performance. Therefore, through this work,  
301 we use domain knowledge to suppress erroneous actions when available data is limited. However, this  
302 inclusion may facilitate harmful behavior in the presence of biased domain knowledge. Therefore,  
303 we advocate the use of well-regulated domain knowledge obtained from experts. Beyond this, we do  
304 not foresee any ethical impact on our work.

## 305 References

- 306 [1] A survey on offline reinforcement learning: Taxonomy, review, and open problems. *IEEE Transactions on*  
307 *Neural Networks and Learning Systems*, 2023. ISSN 21622388. doi: 10.1109/TNNLS.2023.3250269.
- 308 [2] Rishabh Agarwal, Dale Schuurmans, and Mohammad Norouzi. An optimistic perspective on offline  
309 reinforcement learning. In *International Conference on Machine Learning*, pages 104–114. PMLR, 2020.
- 310 [3] Gaon An, Seungyong Moon, Jang-Hyun Kim, and Hyun Oh Song. Uncertainty-based offline reinforcement  
311 learning with diversified q-ensemble. *Advances in neural information processing systems*, 34:7436–7447,  
312 2021.
- 313 [4] Peter L. Bartlett and Ambuj Tewari. Regal: A regularization based algorithm for reinforcement learning in  
314 weakly communicating mdps. In *Proceedings of the Twenty-Fifth Conference on Uncertainty in Artificial*  
315 *Intelligence*, UAI '09, page 35–42, Arlington, Virginia, USA, 2009. AUAI Press. ISBN 9780974903958.
- 316 [5] Greg Brockman, Vicki Cheung, Ludwig Pettersson, Jonas Schneider, John Schulman, Jie Tang, and  
317 Wojciech Zaremba. Openai gym. *arXiv preprint arXiv:1606.01540*, 2016.
- 318 [6] Maxime Chevalier-Boisvert, Bolun Dai, Mark Towers, Rodrigo de Lazcano, Lucas Willems, Salem Lahlou,  
319 Suman Pal, Pablo Samuel Castro, and Jordan Terry. Minigrid & miniworld: Modular & customizable  
320 reinforcement learning environments for goal-oriented tasks. *CoRR*, abs/2306.13831, 2023.
- 321 [7] Will Dabney, Mark Rowland, Marc Bellemare, and Rémi Munos. Distributional reinforcement learning  
322 with quantile regression. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 32,  
323 2018.
- 324 [8] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. Bert: Pre-training of deep bidirec-  
325 tional transformers for language understanding. *arXiv preprint arXiv:1810.04805*, 2018.
- 326 [9] Damien Ernst, Pierre Geurts, and Louis Wehenkel. Tree-based batch mode reinforcement learning. *Journal*  
327 *of Machine Learning Research*, 6, 2005.
- 328 [10] Justin Fu, Aviral Kumar, Ofir Nachum, George Tucker, and Sergey Levine. D4rl: Datasets for deep  
329 data-driven reinforcement learning. *arXiv preprint arXiv:2004.07219*, 2020.
- 330 [11] Scott Fujimoto and Shixiang Shane Gu. A minimalist approach to offline reinforcement learning. *Advances*  
331 *in neural information processing systems*, 34:20132–20145, 2021.
- 332 [12] Scott Fujimoto, Edoardo Conti, Mohammad Ghavamzadeh, and Joelle Pineau. Benchmarking batch deep  
333 reinforcement learning algorithms. *arXiv preprint arXiv:1910.01708*, 2019.
- 334 [13] Scott Fujimoto, David Meger, and Doina Precup. Off-policy deep reinforcement learning without explo-  
335 ration. In *International conference on machine learning*, pages 2052–2062. PMLR, 2019.
- 336 [14] Yarın Gal and Zoubin Ghahramani. A theoretically grounded application of dropout in recurrent neural  
337 networks. *Advances in neural information processing systems*, 29, 2016.
- 338 [15] Sinong Geng, Aldo Pacchiano, Andrey Kolobov, and Ching-An Cheng. Improving offline rl by blending  
339 heuristics. *arXiv preprint arXiv:2306.00321*, 2023.
- 340 [16] Caglar Gulcehre, Sergio Gómez Colmenarejo, Ziyu Wang, Jakub Sygnowski, Thomas Paine, Konrad  
341 Zolna, Yutian Chen, Matthew Hoffman, Razvan Pascanu, and Nando de Freitas. Regularized behavior  
342 value estimation. *arXiv preprint arXiv:2103.09575*, 2021.
- 343 [17] Geoffrey Hinton, Oriol Vinyals, and Jeff Dean. Distilling the knowledge in a neural network. *arXiv*  
344 *preprint arXiv:1503.02531*, 2015.
- 345 [18] Zhiting Hu, Xuezhe Ma, Zhengzhong Liu, Eduard Hovy, and Eric Xing. Harnessing deep neural networks  
346 with logic rules. *arXiv preprint arXiv:1603.06318*, 2016.

- 347 [19] Sham Kakade and John Langford. Approximately optimal approximate reinforcement learning. In  
348 *Proceedings of the Nineteenth International Conference on Machine Learning*, pages 267–274, 2002.
- 349 [20] Ilya Kostrikov, Rob Fergus, Jonathan Tompson, and Ofir Nachum. Offline reinforcement learning with  
350 fisher divergence critic regularization. In *International Conference on Machine Learning*, pages 5774–5783.  
351 PMLR, 2021.
- 352 [21] Aviral Kumar, Justin Fu, Matthew Soh, George Tucker, and Sergey Levine. Stabilizing off-policy q-learning  
353 via bootstrapping error reduction. *Advances in Neural Information Processing Systems*, 32, 2019.
- 354 [22] Aviral Kumar, Aurick Zhou, George Tucker, and Sergey Levine. Conservative q-learning for offline  
355 reinforcement learning. *Advances in Neural Information Processing Systems*, 33:1179–1191, 2020.
- 356 [23] Sergey Levine, Aviral Kumar, G. Tucker, and Justin Fu. Offline reinforcement learning: Tutorial,  
357 review, and perspectives on open problems. *ArXiv*, abs/2005.01643, 2020. URL <https://api.semanticscholar.org/CorpusID:218486979>.
- 358 [24] Sergey Levine, Aviral Kumar, George Tucker, and Justin Fu. Offline reinforcement learning: Tutorial,  
359 review, and perspectives on open problems. *arXiv preprint arXiv:2005.01643*, 2020.
- 361 [25] Siqi Liu, Kay Choong See, Kee Yuan Ngiam, Leo Anthony Celi, Xingzhi Sun, and Mengling Feng.  
362 Reinforcement learning for clinical decision support in critical care: comprehensive review. *Journal of  
363 medical Internet research*, 22(7):e18477, 2020.
- 364 [26] Volodymyr Mnih, Koray Kavukcuoglu, David Silver, Andrei A Rusu, Joel Veness, Marc G Bellemare,  
365 Alex Graves, Martin Riedmiller, Andreas K Fiedjeland, Georg Ostrovski, et al. Human-level control through  
366 deep reinforcement learning. *nature*, 518(7540):529–533, 2015.
- 367 [27] Andrew William Moore. Efficient memory-based learning for robot control. Technical report, University  
368 of Cambridge, Computer Laboratory, 1990.
- 369 [28] Susan A Murphy. A generalization error for q-learning. 2005.
- 370 [29] Dean A Pomerleau. Efficient training of artificial neural networks for autonomous navigation. *Neural  
371 computation*, 3(1):88–97, 1991.
- 372 [30] Rong-Jun Qin, Xingyuan Zhang, Songyi Gao, Xiong-Hui Chen, Zewen Li, Weinan Zhang, and Yang Yu.  
373 Neorl: A near real-world benchmark for offline reinforcement learning. *Advances in Neural Information  
374 Processing Systems*, 35:24753–24765, 2022.
- 375 [31] Kajetan Schweighofer, Markus Hofmarcher, Marius-Constantin Dinu, Philipp Renz, Angela Bitto-Nemling,  
376 Vihang Prakash Patil, and Sepp Hochreiter. Understanding the effects of dataset characteristics on offline  
377 reinforcement learning. In *Deep RL Workshop NeurIPS 2021*, 2021. URL <https://openreview.net/forum?id=A4EWtf-T03Y>.
- 378 [32] Kajetan Schweighofer, Marius-constantin Dinu, Andreas Radler, Markus Hofmarcher, Vihang Prakash  
379 Patil, Angela Bitto-Nemling, Hamid Eghbal-zadeh, and Sepp Hochreiter. A dataset perspective on offline  
380 reinforcement learning. In *Conference on Lifelong Learning Agents*, pages 470–517. PMLR, 2022.
- 381 [33] Andrew Silva and Matthew Gombolay. Encoding human domain knowledge to warm start reinforcement  
382 learning. In *Proceedings of the AAAI conference on artificial intelligence*, volume 35, pages 5042–5050,  
383 2021.
- 384 [34] Andrew Silva, Matthew Gombolay, Taylor Killian, Ivan Jimenez, and Sung-Hyun Son. Optimization  
385 methods for interpretable differentiable decision trees applied to reinforcement learning. In *International  
386 conference on artificial intelligence and statistics*, pages 1855–1865. PMLR, 2020.
- 387 [35] Samarth Sinha, Ajay Mandlekar, and Animesh Garg. S4rl: Surprisingly simple self-supervision for offline  
388 reinforcement learning in robotics. In Aleksandra Faust, David Hsu, and Gerhard Neumann, editors,  
389 *Proceedings of the 5th Conference on Robot Learning*, volume 164 of *Proceedings of Machine Learning  
390 Research*, pages 907–917. PMLR, 08–11 Nov 2022.
- 391 [36] Kihyuk Sohn, Zizhao Zhang, Chun-Liang Li, Han Zhang, Chen-Yu Lee, and Tomas Pfister. A simple  
392 semi-supervised learning framework for object detection. *arXiv preprint arXiv:2005.04757*, 2020.
- 393 [37] Jiayi Tang and Ke Wang. Ranking distillation: Learning compact ranking models with high performance  
394 for recommender system. In *Proceedings of the 24th ACM SIGKDD international conference on knowledge  
395 discovery & data mining*, pages 2289–2298, 2018.
- 396

- 397 [38] Raphael Tang, Yao Lu, Linqing Liu, Lili Mou, Olga Vechtomova, and Jimmy Lin. Distilling task-specific  
398 knowledge from bert into simple neural networks. *arXiv preprint arXiv:1903.12136*, 2019.
- 399 [39] Wei-Cheng Tseng, Tsun-Hsuan Johnson Wang, Yen-Chen Lin, and Phillip Isola. Offline multi-agent  
400 reinforcement learning with knowledge distillation. In S. Koyejo, S. Mohamed, A. Agarwal, D. Belgrave,  
401 K. Cho, and A. Oh, editors, *Advances in Neural Information Processing Systems*, volume 35, pages  
402 226–237. Curran Associates, Inc., 2022.
- 403 [40] Richa Verma, Durgesh Kalwar, Harshad Khadilkar, and Balaraman Ravindran. Guiding offline reinforce-  
404 ment learning using a safety expert. In *Proceedings of the 7th Joint International Conference on Data  
405 Science & Management of Data (11th ACM IKDD CODS and 29th COMAD)*, pages 82–90, 2024.
- 406 [41] Ziyu Wang, Alexander Novikov, Konrad Zolna, Josh S Merel, Jost Tobias Springenberg, Scott E Reed,  
407 Bobak Shahriari, Noah Siegel, Caglar Gulcehre, Nicolas Heess, et al. Critic regularized regression.  
408 *Advances in Neural Information Processing Systems*, 33:7768–7778, 2020.
- 409 [42] Yue Wu, Shuangfei Zhai, Nitish Srivastava, Joshua M. Susskind, Jian Zhang, Ruslan Salakhutdinov, and  
410 Hanlin Goh. Uncertainty weighted actor-critic for offline reinforcement learning. In *International Confer-  
411 ence on Machine Learning*, 2021. URL <https://api.semanticscholar.org/CorpusID:234763307>.
- 412 [43] Qizhe Xie, Minh-Thang Luong, Eduard Hovy, and Quoc V Le. Self-training with noisy student improves  
413 imagenet classification. In *Proceedings of the IEEE/CVF conference on computer vision and pattern  
414 recognition*, pages 10687–10698, 2020.
- 415 [44] Qisen Yang, Shenzhi Wang, Matthieu Gaetan Lin, Shiji Song, and Gao Huang. Boosting offline reinforce-  
416 ment learning with action preference query. In *International Conference on Machine Learning*, pages  
417 39509–39523. PMLR, 2023.
- 418 [45] Li Yuan, Francis EH Tay, Guilin Li, Tao Wang, and Jiashi Feng. Revisiting knowledge distillation via label  
419 smoothing regularization. In *2020 IEEE/CVF Conference on Computer Vision and Pattern Recognition  
420 (CVPR)*, pages 3902–3910, 2020. doi: 10.1109/CVPR42600.2020.00396.
- 421 [46] Xiaoxuan Zhang and S Zhang Y Yu. Domain knowledge guided offline q learning. In *Second Offline  
422 Reinforcement Learning Workshop at Neurips*, volume 2021, 2021.
- 423 [47] Ying Zheng, Haoyu Chen, Qingyang Duan, Lixiang Lin, Yiyang Shao, Wei Wang, Xin Wang, and  
424 Yuedong Xu. Leveraging domain knowledge for robust deep reinforcement learning in networking.  
425 In *IEEE INFOCOM 2021 - IEEE Conference on Computer Communications*, pages 1–10, 2021. doi:  
426 10.1109/INFOCOM42981.2021.9488863.

427 **A Theoretical Analysis**

428 **Notations**

429 For any deterministic policy  $\pi$  the performance return is formulated as  $\eta(\pi) = \mathbb{E}_{\tau \sim \pi} [\sum_{t=0}^{\infty} \gamma^t r(s_t, a_t)]$

430 For any policy  $\pi$ ,  $\rho_\pi$  is the (unnormalized) discounted visitation frequency given by  $\rho_\pi(s) = \sum_{t=0}^{\infty} \gamma^t P(s_t = s)$   
 431 where  $s_0 \sim \rho^0(s_0)$  and the trajectory  $(s_0, s_1, \dots)$  is sampled from the policy  $\pi$  and  $\rho_\pi(s) \in [0, \frac{1}{1-\gamma}]$ .

432  $\bar{\rho}_\pi(s) = \sup\{\rho_\pi(s), s \in S\} \in [\frac{1}{|S_\pi|(1-\gamma)}, \frac{1}{(1-\gamma)}]$

433 We denote the regularized policy learned by ExID on  $\mathcal{B}_r$  as  $\hat{\pi}$  and the unregularized policy as  $\pi_u$ .

434 **Lemmas**

435 We introduce the following Lemma required for our theoretical analysis.

436 **Lemma A.1.** ([44]) Given two policies  $\pi_1$  and  $\pi_2$

$$\eta(\pi_1) - \eta(\pi_2) = \int_{s \in S} \rho_{\pi_1}(s)(Q^*(s, \pi_1(s)) - V^*(s))ds - \int_{s \in S} \rho_{\pi_2}(s)(Q^*(s, \pi_2(s)) - V^*(s))ds$$

437 *Proof.* Please refer to Lemma A.1 Eq 17 in [44] □

438 **Proposition A.2.** (4.2) Denote  $\hat{\pi}$  as the policy learned by ExID,  $\pi_u$  as any offline RL policy learned on  $\mathcal{B}_r$  and  
 439 optimal  $Q$  function as  $Q^*$  and  $V$  function as  $V^*$ . Then it holds that

$$\eta(\hat{\pi}) - \eta(\pi_u) \geq \mathbb{E}_{s \sim O | \pi_u} [V^*(s) - Q^*(s, \pi_u(s))] - \bar{\rho}_{\hat{\pi}} \alpha$$

440 *Proof.* According to [19] performance improvement between two policies if given by

$$\eta(\pi_1) = \eta(\pi_2) + \mathbb{E}_{\tau \sim \pi_1} \left[ \sum_{t=0}^{\infty} \gamma^t Q_{\pi_2}(s_t, a_t) - V_{\pi_2}(s_t) \right] \quad (8)$$

441 Replacing  $\pi_1$  by  $\hat{\pi}$  and  $\pi_2$  by  $\pi_u$  and by following Lemma A.1

$$\eta(\hat{\pi}) - \eta(\pi_u) = \int_{s \in S} \rho_{\hat{\pi}}(s)(Q^*(s, \hat{\pi}(s)) - V^*(s))ds - \int_{s \in S} \rho_{\pi_u}(s)(Q^*(s, \pi_u(s)) - V^*(s))ds \quad (9)$$

$$= \int_{s \in S} \rho_{\pi_u}(s)(V^*(s) - Q^*(s, \pi_u(s)))ds - \int_{s \in S} \rho_{\hat{\pi}}(s)(V^*(s) - Q^*(s, \hat{\pi}(s)))ds \quad (10)$$

442 Dividing the state space into in dataset domain states (I) and OOD states (O). The

$$\underbrace{\left[ \int_{s \in I} \rho_{\pi_u}(s)(V^*(s) - Q^*(s, \pi_u(s)))ds - \int_{s \in I} \rho_{\hat{\pi}}(s)(V^*(s) - Q^*(s, \hat{\pi}(s)))ds \right]}_a + \underbrace{\left[ \int_{s \in O} \rho_{\pi_u}(s)(V^*(s) - Q^*(s, \pi_u(s)))ds - \int_{s \in O} \rho_{\hat{\pi}}(s)(V^*(s) - Q^*(s, \hat{\pi}(s)))ds \right]}_b \quad (11)$$

Since the regularization loss facilitates visitation to OOD states via knowledge distillation we assume

443  $\rho_{\hat{\pi}} = \rho_{\pi_u} - \Delta_i$  for  $s \in i$  and  $\rho_{\hat{\pi}} = \rho_{\pi_u} + \Delta_o$  for  $s \in o$  where  $\Delta_i \in [0, \rho_{\pi_u}(s)]$  and  $\Delta_o \in [0, \frac{1}{1-\gamma} - \rho_{\pi_u}(s)]$

$$a = \int_{s \in I} \rho_{\pi_u}(s)(V^*(s) - Q^*(s, \pi_u(s)))ds - \int_{s \in I} (\rho_{\pi_u} - \Delta_i)(s)(V^*(s) - Q^*(s, \hat{\pi}(s)))ds \quad (13)$$

$$= \int_{s \in I} \rho_{\pi_u}(s)(Q^*(s, \hat{\pi}(s)) - Q^*(s, \pi_u(s)))ds + \int_{s \in I} \Delta_i(s)(V^*(s) - Q^*(s, \hat{\pi}(s)))ds \quad (14)$$

Under assumption in distribution action can be learned from the dataset due to conservatism of offline RL

$$444 \quad (Q^*(s, \hat{\pi}(s)) - Q^*(s, \pi_u(s))) \approx 0, a \geq 0$$

$$b = \int_{s \in O} \rho_{\pi_u}(s)(V^*(s) - Q^*(s, \pi_u(s)))ds - \int_{s \in O} (\rho_{\pi_u} + \Delta_o)(s)(V^*(s) - Q^*(s, \hat{\pi}(s)))ds \quad (15)$$

$$\geq \int_{s \in O} \rho_{\pi_u}(s)(V^*(s) - Q^*(s, \pi_u(s)))ds - \int_{s \in O} \rho_{\hat{\pi}}(s)(V^*(s) - Q^*(s, \hat{\pi}(s)))ds \quad (16)$$

$$\geq \mathbb{E}_{s \sim O | \pi_u}[V^*(s) - Q^*(s, \pi_u(s))] - \mathbb{E}_{s \sim O | \hat{\pi}}[V^*(s) - Q^*(s, \hat{\pi}(s))] \quad (17)$$

445 Further loosening the lower bound

$$= \mathbb{E}_{s \sim O | \pi_u}[V^*(s) - Q^*(s, \pi_u(s))] - \bar{\rho}_{\hat{\pi}} \int_{s \in O} \frac{\rho_{\hat{\pi}}}{\bar{\rho}_{\hat{\pi}}}(V^*(s) - Q^*(s, \hat{\pi}(s)))ds \quad (18)$$

$$\geq \mathbb{E}_{s \sim O | \pi_u}[V^*(s) - Q^*(s, \pi_u(s))] - \bar{\rho}_{\hat{\pi}} \int_{s \in O} (V^*(s) - Q^*(s, \hat{\pi}(s)))ds \quad (19)$$

446 Combining Eq 14, 17 and 19, and denoting  $\alpha = \mathbb{E}_{s \sim O}[V^*(s) - Q^*(s, \hat{\pi}(s))]$

$$\eta(\hat{\pi}) - \eta(\pi_u) \geq \mathbb{E}_{s \sim O | \pi_u}[V^*(s) - Q^*(s, \pi_u(s))] - \bar{\rho}_{\hat{\pi}} \alpha \quad (20)$$

447 Hence, Proposition 4.2 follows **Q.E.D**

448

□

449 **Proposition A.3.** (4.3) Algo 1 reduces generalization error if  $Q^*(s, \pi_t^\omega(s)) > Q^*(s, \pi(s))$  for  $s \in \mathcal{D} \cap \mathcal{B}_r$ ,  
450 where  $\pi$  is vanilla offline RL policy learnt on  $\mathcal{B}_r$ .

451 *Proof.* Generalization error for any policy  $\pi$  as defined by [28] can be written as:

$$G_\pi = V^*(s_0) - V_\pi(s_0) = -\mathbb{E}_{\tau \sim \pi} \left[ \sum_{t=0}^T \gamma^t Q^*(s_t, \pi(s_t)) - V^*(s_t) \right] \quad (21)$$

452 Here,  $\mathbb{E}_{\tau \sim \pi}$  represents sampling trajectories with policy  $\pi$ . Since the state space is continuous, we can represent the expectation as an integral over the state space

$$= - \sum_{t=0}^T \gamma^t \int_{s \in S} P(s_t = s | \pi) (Q^*(s_t, \pi(s_t)) - V^*(s_t)) ds \quad (22)$$

$$= - \int_{s \in S} \sum_{t=0}^T \gamma^t P(s_t = s | \pi) (Q^*(s_t, \pi(s_t)) - V^*(s_t)) ds \quad (23)$$

453 Analysing with respect to  $s \in \mathcal{D} \cap \mathcal{B}_r$  we can break the integration into two parts

$$= - \left[ \int_{s \in S/D} \sum_{t=0}^T \gamma^t P(s_t = s | \pi) (Q^*(s_t, \pi(s_t)) - V^*(s_t)) ds + \int_{s \in D} \sum_{t=0}^T \gamma^t P(s_t = s | \pi) (Q^*(s_t, \pi(s_t)) - V^*(s_t)) \right] \quad (24)$$

$$= - \left[ f(s | \pi) + \int_{s \in D} \sum_{t=0}^T \gamma^t P(s_t = s | \pi) (Q^*(s_t, \pi(s_t)) - V^*(s_t)) \right] \quad (25)$$

For a policy  $\hat{\pi}$  learnt in Algo 1 the action for  $s_t = s \in \mathcal{D}$  is regularized to be close to  $\pi_t^\omega$  which either follows domain knowledge or expert demonstrations. Hence, it is reasonable to assume  $Q^*(s_t, \pi_t^\omega(s_t)) > Q^*(s_t, \pi(s_t))$ .

454 It follows

$$\int_{s \in D} \sum_{t=0}^T \gamma^t P(s_t = s | \hat{\pi}) (Q^*(s_t, \hat{\pi}(s_t)) - V^*(s_t)) < \int_{s \in D} \sum_{t=0}^T \gamma^t P(s_t = s | \pi) (Q^*(s_t, \pi(s_t)) - V^*(s_t)) \quad (26)$$

455 **Note for**  $s \notin \mathcal{D}$ ,  $f(s|\hat{\pi}) \approx f(s|\pi)$ . This is because the regularization term assigns max Q value to a different action for  $s \in \mathcal{D}$  but  $\max_a(Q(s, a))$  remains same

$$\begin{aligned} & \therefore - \left[ f(s|\hat{\pi}) + \int_{s \in \mathcal{D}} \sum_{t=0}^T \gamma^t P(s_t = s|\hat{\pi})(Q^*(s_t, \hat{\pi}(s_t)) - V^*(s_t)) \right] \\ & < - \left[ f(s|\pi) + \int_{s \in \mathcal{D}} \sum_{t=0}^T \gamma^t P(s_t = s|\pi)(Q^*(s_t, \pi(s_t)) - V^*(s_t)) \right] \end{aligned} \quad (27)$$

456 Hence,  $G_{\hat{\pi}} < G_{\pi}$  Proposition 2 follows **Q.E.D**

457

□

458 **B Missing Examples**

459 *Performing Q – Learning by sampling from a reduced batch  $\mathcal{B}_r$  may not converge to an optimal policy for the*  
 460 *MDP  $M_{\mathcal{B}}$  representing the full buffer.*

461 **Example** (Theorem 1,[13]) defines MDP  $M_{\mathcal{B}}$  of  $\mathcal{B}$  from same state action space of the original MDP  $M$  with  
 462 transition probabilities  $p_{\mathcal{B}}(s'|s, a) = \frac{N(s, a, s')}{\sum_{\bar{s}} N(s, a, \bar{s})}$  where  $N(s, a, s')$  is the number of times  $(s, a, s')$  occurs in  $\mathcal{B}$   
 463 and an terminal state  $s_{init}$ . It states  $p_{\mathcal{B}}(s_{init}|s, a) = 1$  when  $\sum_{\bar{s}} N(s, a, \bar{s}) = 0$ . This happens when transitions  
 464 of some  $s'$  of  $(s, a, s')$  are missing from the buffer, which may occur in  $\mathcal{B}_r$  when  $\mathcal{B}_r \subset \mathcal{B}$ .  $r(s_{init}, s, a)$  is  
 465 initialized to  $Q(s, a)$ . We assume that a policy learned on reduced dataset  $\mathcal{B}_r$  converges to optimal value function  
 466 and disprove it using the following counterexample:

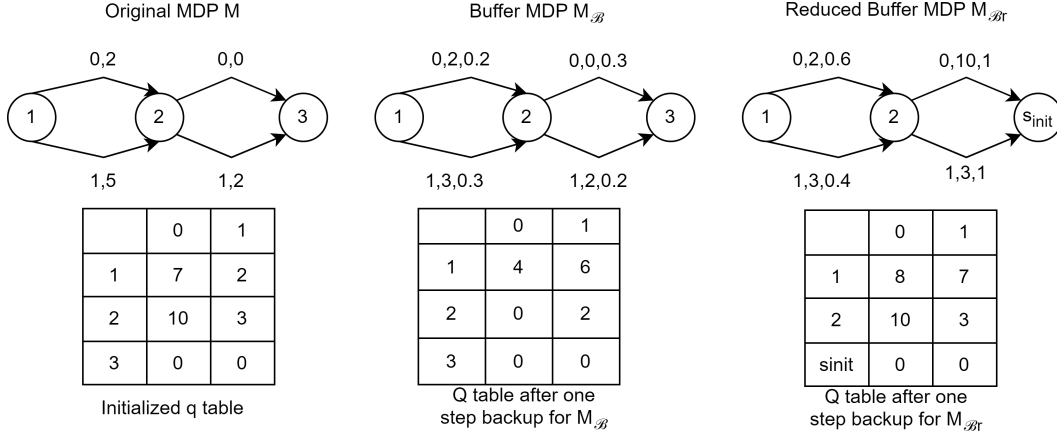
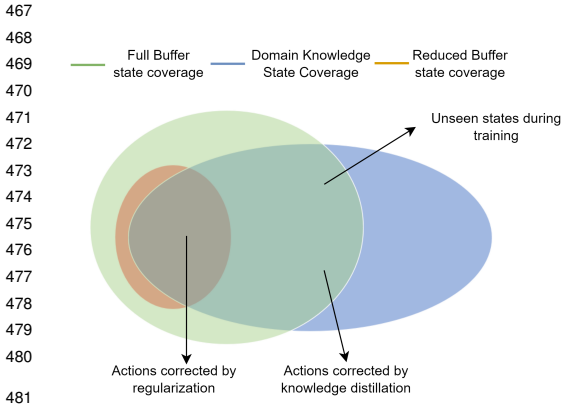


Figure 7: Example MDP, sampled buffer MDP and reduced buffer with Q tables



482 **Figure 8:** We hypothesize the suboptimal performance of offline RL for limited data can be addressed via domain knowledge via action regularization and knowledge distillation.

486

487

488

489 A visualization is shown in Fig 8.

490 **C Algorithm**

491 The pseudo code of the algorithm is described in Algo 1.

We take a simple MDP illustrated in Fig 7 with 3 states and 2 actions (0,1). The reward of each action is marked along the transition. The sampled MDP is constructed the following samples (1,0,2)-2,(1,1,2)-3, (2,0,3)-3, and (2,1,3)-2 and the reduced buffer MDP with samples (1,0,2)-2 and (1,1,2)-1. The probabilities are marked along the transition. It is easy to see that the policy learned under the reduced MDP converges to a nonoptimal policy after one step of the Q table update with  $Q(s, a) = r(s, a) + p(s'|s, a) * \max_{a'} (Q(s', a'))$ . This happens because of transition probability shift on reducing samples  $p_{\mathcal{B}}(s'|s, a) \neq p_{\mathcal{B}_r}(s'|s, a)$  and no Q updates for  $(s, a) \notin \mathcal{B}_r$ .

Our methodology addresses these issues as follows:

- For  $s \in D \cap \mathcal{B}_r$ , better actions are enforced through regularization using  $\pi_i^\omega$  even when the transition probabilities are low for optimal transitions.
- Incorporating regularization distills the teacher’s knowledge in the critic-enhancing generalization.



---

**Algorithm 1** Pseudo code for EXID

---

```
1: Input: Reduced buffer  $\mathcal{B}_r$ , Initial teacher network  $\pi_t^\omega$ , Training steps  $N$ , Warm-up steps  $k$ , Soft
   update  $\tau$ , hyperparameters:  $\lambda, \alpha$ 
2: Initialize Critic with MC dropout and Target Critic  $Q_s^\theta, Q_s^{\theta'}$ 
3: for  $n \leftarrow 1$  to  $N$  do
4:   Sample mini-batch  $b$  of transitions  $(s, a, r, s') \sim \mathcal{B}_r$   $a_t = [], a_s = [], s_r = []$ 
5:   for  $s \in b$  do
6:     if  $s \models \mathcal{D}$  and  $\pi_t^\omega(s) \neq \operatorname{argmax}_a(Q_s^\theta(s, a))$  then
7:        $a_t.append(\pi_t^\omega(s))$ 
8:        $a_s.append(\operatorname{argmax}_a(Q_s^\theta(s, a)))$ 
9:        $s_r.append(s)$ 
10:    end if
11:  end for
12:  if  $n > k \wedge \text{Cond. 6}$  then
13:    Update  $\pi_t^\omega(s)$  using Eq 7
14:     $\mathcal{L}_r(\theta) = 0$ 
15:  else
16:    Calculate  $\mathcal{L}_r(\theta)$  using Eq 3
17:  end if
18:  Calculate  $\mathcal{L}(\theta)$  using Eq 4
19:  Update  $Q_s^\theta$  with  $\mathcal{L}(\theta)$  and softy update  $Q_s^{\theta'}$  and  $\tau$ 
20: end for
```

---

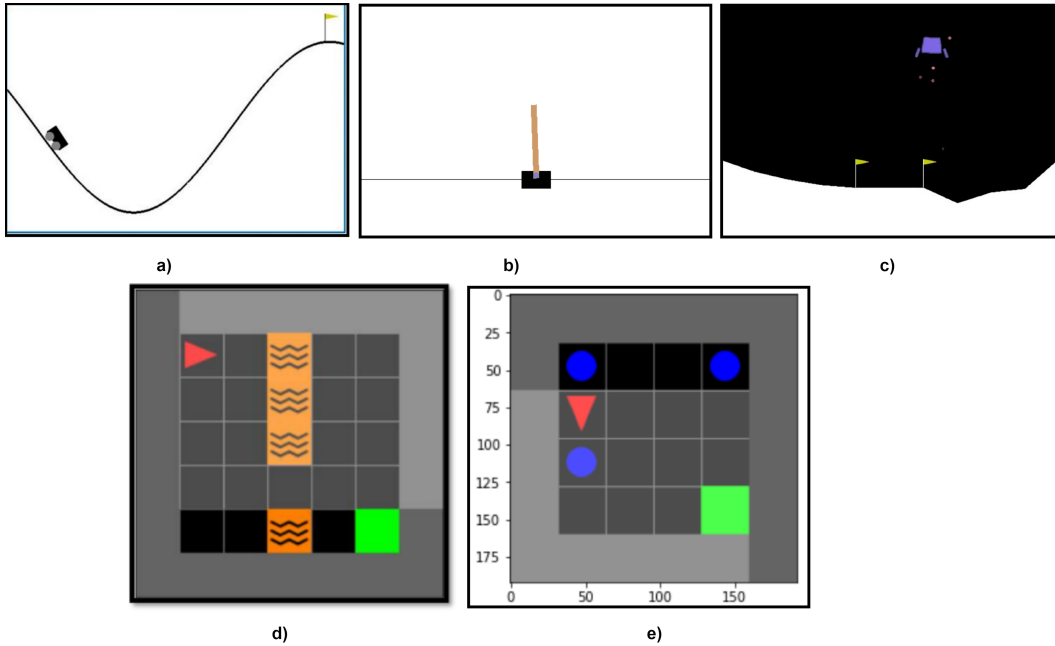
492 **D Environments and Domain Knowledge Trees**

Figure 9: Graphical visualizations of environments used in the experiments. These environments are a) MountainCar-v0 b) CartPole-v1 c) LunarLander-v2 d) MiniGrid-LavaGapS7-v0 e) MiniGrid-Dynamic-Obstacles-Random-6x6-v0

493 The graphical visualization of each environment is depicted in Fig 9. The choice of environment in this paper  
494 depended on two factors: a) Pre-existing standard methods of generating offline RL datasets. b) Possibility of  
495 creating intuitive decision tree-based domain knowledge. All datasets have been created via [31]. We explain the  
496 environments in detail as follows:

497 **Mountain-car Environment:** This environment Fig 9 a) has two state variables, position and velocity, and three  
 498 discrete actions: left push, right push, and no action [27]. The goal is to drive a car up a valley to reach the flag.  
 499 This environment is challenging for offline RL because of sparse rewards, which are only obtained on reaching  
 500 the flag.

501 **Cart-pole Environment** The environment Fig 9 b) has 4 states and 2 actions representing left force and right  
 502 force. The objective is to balance a pole on a moving cart.

503 **Lunar-Lander Environment:** The task is to land a lunar rover between two flags Fig 9 c) by observing 8 states  
 504 and applying one of 4 actions.

505 **Minigrd Environments:** Mini-grid [6] is an environment suite containing 2D grid-worlds with goal oriented  
 506 tasks. As explained in the main text, we experiment using MiniGrid-LavaGapS7-v0 and MiniGrid-Dynamic-  
 507 Obstacles-Random-6x6-v0 from this environment suite is shown in Fig 9 d) and e). In MiniGrid-LavaGapS7-v0,  
 508 the agent has to avoid Lava and pass through the gap to reach the goal. Dynamic obstacles are similar; however,  
 509 the agent can start at a random position and has to avoid dynamically moving balls to reach the goal. The  
 510 environment has image observation with 3 channels (OBJECT\_ID, COLOR\_ID, STATE). Following [31]  
 511 experiments, we flatten the image to an array of 98 observations and restrict action space to three actions: Turn  
 512 left, Turn Right, and Move forward. The results of minigrd environment are reported in Table 3. Since this  
 513 environment uses a semantic map from image observation, we collect states from a fixed policy with random  
 514 actions to generate the teacher’s state distribution. CQL on the full dataset achieves the average reward of  
 515  $0.92 \pm 0.1$  for DynamicObstacles and  $0.53 \pm 0.01$  for LavaGapS.

516 The domain knowledge trees for all the environments are shown in Fig 10. The cart pole domain knowledge  
 517 tree Fig 10 a) is taken from [33] (Fig 7). The Lunar Lander decision nodes Fig 10 b) have been taken from [34]  
 518 (Fig4). For the mini-grid environments, we construct intuitive decision trees shown in Fig 10 d) and Fig 10 e).  
 519 Positions 52, 40, and 68 represent positions front, right, and left of the agent. Value 0.2 represents a wall, 0.9  
 520 represents Lava, and 0.6 represents a ball. We check positions 52, 40, and 68 for these obstacles and choose the  
 521 recommended actions as domain knowledge.

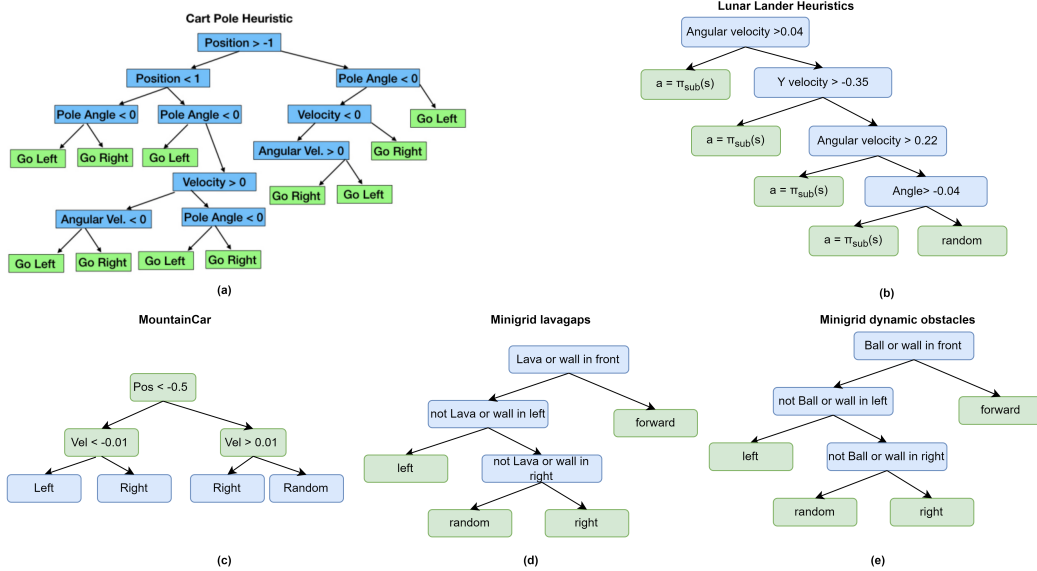


Figure 10: Domain knowledge trees for a) CartPole-v1 b) LunarLander-v2 c) MountainCar-v0 d) MiniGrid-LavaGapS7-v0 e) MiniGrid-Dynamic-Obstacles-Random-6x6-v0 environments

## 522 E Related Work: Knowledge Distillation

523 Knowledge distillation is a well-embraced technique of incorporating additional information in neural networks  
 524 and has been applied to various fields like computer vision [43, 36], natural language processing [8, 38], and  
 525 recommendation systems [37]. [17] introduced the concept of distilling knowledge from a complex, pre-trained  
 526 model (teacher) into a smaller model (student). In recent years, researchers have explored the integration  
 527 of rule-based regularization techniques within the context of knowledge distillation. Rule regularization  
 528 introduces additional constraints based on predefined rules, guiding the learning process of the student model

Table 3: Average reward  $\uparrow$  obtained during online evaluation over 3 seeds on Minigrid environments

ENVIRONMENT	$\mathcal{D}$	BC D	BCQ D	CQL D	ExID
MINIGRID	0.50	0.59	0.24	0.14	0.79
DYNAMIC	$\pm$	$\pm$	$\pm$	$\pm$	$\pm$
RANDOM6X6	0.08	0.07	0.22	0.1	0.07
MINIGRID	0.27	0.29	0.26	0.28	0.46
LAVAGAPS	$\pm$	$\pm$	$\pm$	$\pm$	$\pm$
7X7	0.09	0.11	0.1	0.12	0.13

529 [18, 45]. These techniques have shown to reduce overfitting and enhance generalization [38]. Knowledge  
 530 distillation is also prevalent in the field of RL [47] and offline RL [39]. Contrary to prevalent teacher-student  
 531 knowledge distillation techniques, our work does not enforce parameter sharing among the networks. Through  
 532 experiments, we demonstrate that a simple regularization loss and expected performance-based updates can  
 533 improve generalization to unobserved states covered by domain knowledge. There are also no constraints on  
 534 keeping the same network structure for the teacher, paving ways for capturing the domain knowledge into more  
 535 structured networks such as Differentiable Decision Trees (DDTs).

## 536 F Network Architecture and Hyper-parameters

537 We follow the network architecture and hyper-parameters proposed by [31] for all our networks, including the  
 538 baseline networks. The teacher BC network  $\pi_{\omega}^t$  and Critic network  $Q_s^{\theta}(s, a)$  consists of 3 linear layers, each  
 539 having a hidden size of 256 neurons. The number of input and output neurons depends on the environment’s state  
 540 and action size. All layers except the last are SELU activation functions; the final layer uses linear activation.  
 541  $\pi_{\omega}^t$  uses a softmax activation function in the last layer for producing action probabilities. A learning rate of  
 542 0.0001 with batch size 32 and  $\alpha = 0.1$  is used for all environments. MC dropout probability of 0.5 and  
 543 number of stochastic passes  $T=10$  have been used for the critic network. The uncertainty check is performed  
 544 every 15 episodes after the warm start to avoid computational overhead. The hyper-parameters specific to our  
 545 algorithm for OpenAI gym are reported in Table F. The hyper-parameters specific to our algorithm for Minigrid  
 546 environments are reported in Table 5.

Table 4: Hyperparameters for openAI gym environments

HYPERPARAM	MOUNTAINCAR			CARTPOLE			LUNAR-LANDER		
	EXPERT	REPLAY	NOISY	EXPERT	REPLAY	NOISY	EXPERT	REPLAY	NOISY
$\lambda$	0.5	0.5	0.5	0.5	0.5	0.5	0.5	0.5	0.5
$k$	30	30	30	30	30	30	30	30	30
$\pi_{\omega}^t$ LR	$1e^{-5}$	$1e^{-5}$	$1e^{-5}$	$1e^{-2}$	$1e^{-2}$	$1e^{-2}$	$1e^{-4}$	$1e^{-4}$	$1e^{-4}$
TRAINING STEPS	42000	36000	36000	30000	17000	17000	18000	18000	18000

Table 5: Hyper-parameters for Mini-grid environments for replay dataset

Environment	DynamicObstRandom6x6-v0	LavaGapS7v0
$\lambda$	0.1	0.1
$k$	30	30
$\pi_{\omega}^t$ lr	$1e^{-4}$	$1e^{-4}$
training steps	5000	10000

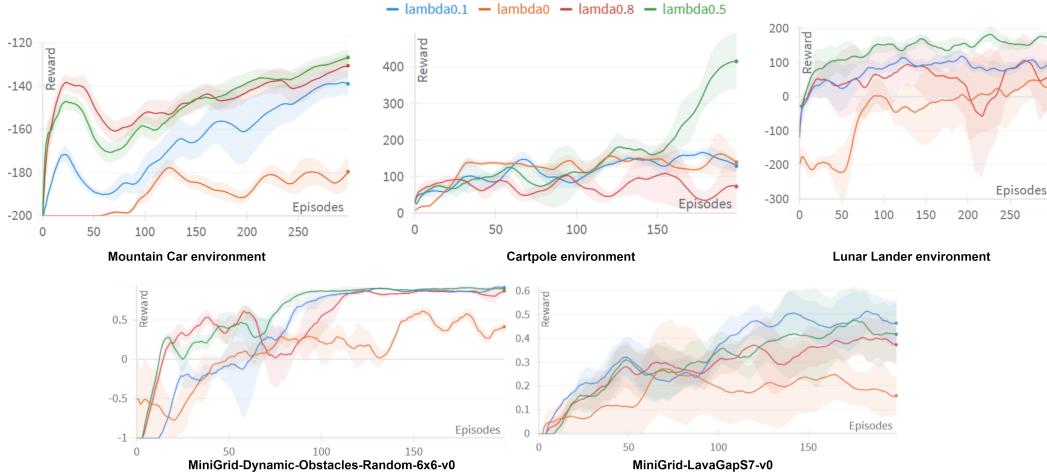


Figure 11: Effect of  $\lambda$  on the performance of ExID for different environments expert datasets.

547 **G Effect of  $k$  and  $\lambda$  and Evaluation Plots**

548 We empirically evaluate the effect of  $\lambda$  In Fig 11 and  $k$  in Fig 12. We believe these parameters depend on the  
 549 quality of  $\mathcal{D}$ . For the given  $\mathcal{D}$  in the environments we empirically observe,  $\lambda = 0.5$  generally performs well,  
 550 except for Minigrid environments where  $\lambda = 0.1$  works better. Increasing the warm start parameter  $k$  generally  
 551 increases the initial performance of the policy, allowing it to learn from the teacher. Meanwhile, no warm start  
 552 adversely affects policy performance as the critic may erroneously update the teacher. From empirical evaluation,  
 553 we observe that  $k = 30$  gives a reasonable start to the policy. All the evaluation plots are shown in Fig 13, where  
 554 it can be observed that ExID performs better than baseline CQL.

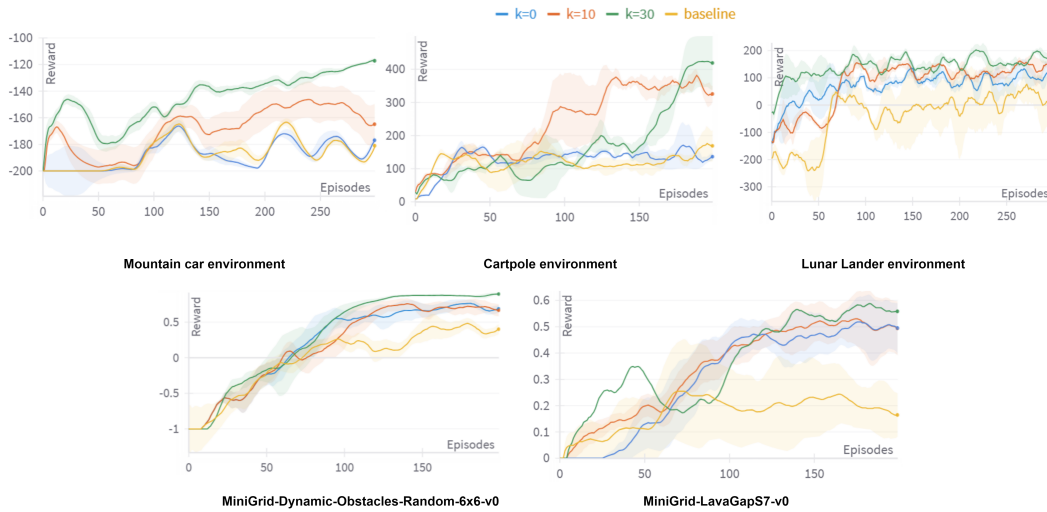


Figure 12: Effect of  $k$  on the performance of ExID for different environments expert datasets.

555 **H Data reduction design and data distribution visualization of reduced**  
 556 **dataset**

557 In this section, we discuss the intuition behind our data-limiting choices. We also visually represent selected  
 558 reduced datasets for the OpenAI gym environments.

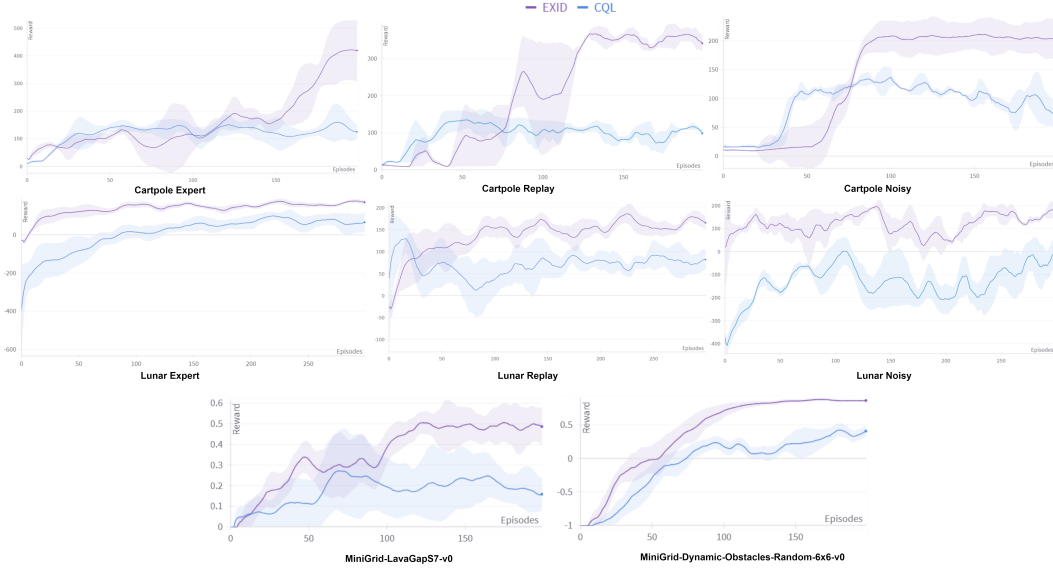


Figure 13: Evaluation plots of CQL and EXID algorithms for Cartpole, Lunar-Lander, and Minigrid environments using different data types and seeds reported in the main paper Table 5.1.

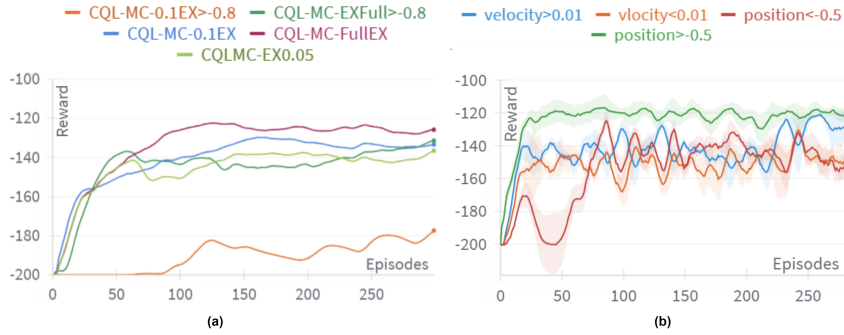


Figure 14: (a) The effect of data reduction and removal on baseline CQL visualized on Mountain Car Environment (b) Performance of ExID on removing different parts of the data based on nodes of Fig 10 (c) from Mountain Car expert dataset

559 **Reducing transitions from the dataset:** For all datasets, 10% of the data samples were extracted from the full  
 560 dataset. This experimental design choice is based on the observation shown in Fig 14 (a). Performance degrades  
 561 on reducing samples to 0.1% of the dataset and reduces further on reducing samples to 0.05% of the dataset.  
 562 However, this drop is not substantial. The performance also reduces on removing part of the dataset from the  
 563 full dataset with states  $> -0.8$ . However, the worst performance is observed when both samples are reduced  
 564 and data is omitted, attributing to accumulated errors from probability ratio shift contributing to an increase in  
 565 generalization error. Our methodology aims to address this gap in performance.

566 **Removing part of the state space:** Due to the simplicity of the Mountain-Car environment, we analyze the  
 567 Mountain-Car expert dataset to show the effect of removing data matching state conditions of the different nodes  
 568 in the decision tree in Fig 10 (c). The performance for each condition is summarised in Table 6. The most  
 569 informative node in the tree is  $\text{position} > -0.5$ ; removing states matching this condition causes a performance  
 570 drop in the algorithm as the domain knowledge regularization does not contribute significant information to the  
 571 policy. Similarly, removing data with  $\text{velocity} < 0.01$  causes a performance drop. However, both performances  
 572 are higher than the baseline CQL trained on reduced data. Based on this observation, we choose state removal  
 573 conditions that preserve states matching part of the information in the tree such that the regularization term  
 574 contributes substantially to the policy. Fig 15 shows the data distribution plot of 10% samples extracted from  
 575 mountain car replay and noisy data with states  $> -0.8$  removed. Fig 16 shows visualizations for 10% samples

576 extracted from expert data with velocity  $> -1.5$  removed. Fig 17 shows visualizations for 10% samples  
 577 extracted from expert data with lander angle  $< -0.04$  removed.

Table 6: Performance of ExID on removing different parts of the data based on nodes of Fig 10 (c) from Mountain Car expert dataset

Position $> -0.5$	Position $< -0.5$	Velocity $> 0.01$	Velocity $< 0.01$
$-121.89 \pm 7.69$	$-151 \pm 13.6$	$-128.48 \pm 11.84$	$-147.80 \pm 5.01$

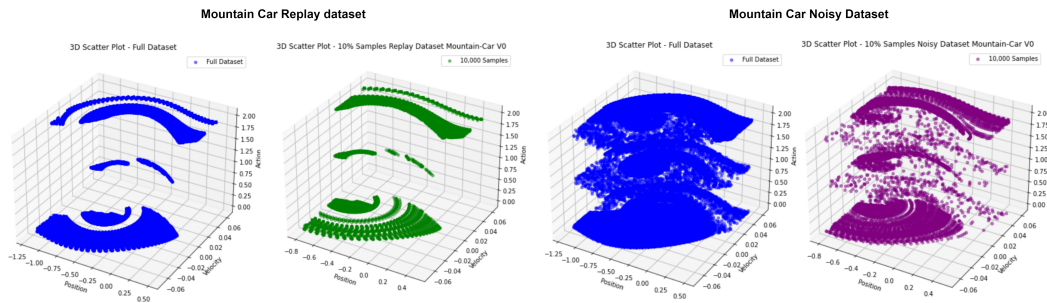


Figure 15: Data distribution of reduced dataset compared to the full dataset for mountain replay and noisy data

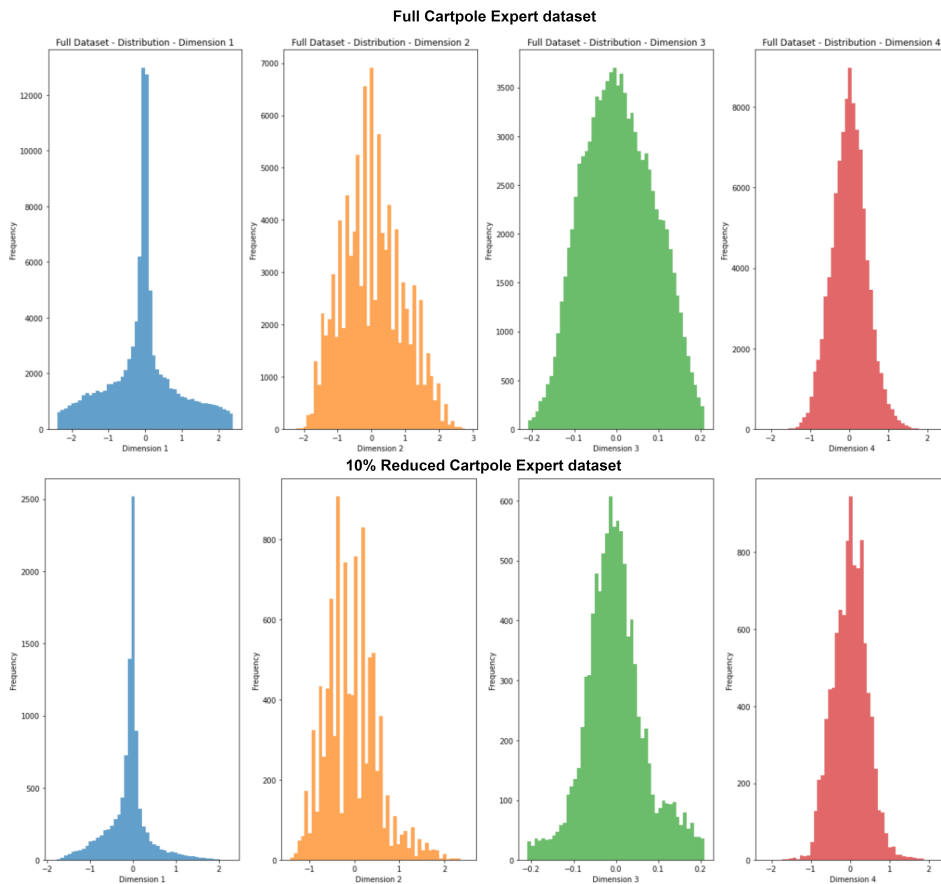


Figure 16: Data distribution of reduced cart pole expert dataset compared to the full dataset

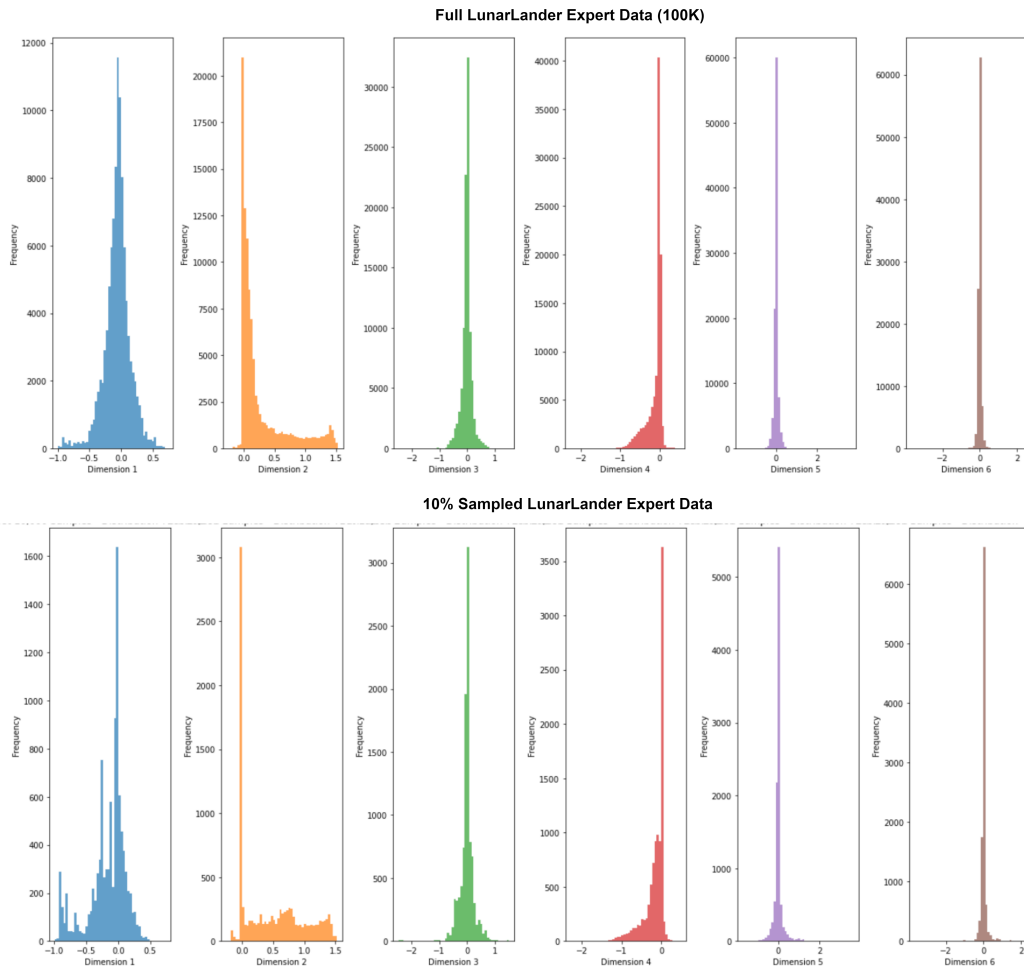


Figure 17: Data distribution of reduced LunarLander expert dataset compared to the full dataset

## 578 **NeurIPS Paper Checklist**

### 579 **1. Claims**

580 Question: Do the main claims made in the abstract and introduction accurately reflect the paper's  
581 contributions and scope?

582 Answer: [\[Yes\]](#)

583 Justification: The claims made in the paper have been experimented on different settings for validity  
584 and generalization. Please refer to sec 5.2.

585 Guidelines:

- 586 • The answer NA means that the abstract and introduction do not include the claims made in the  
587 paper.
- 588 • The abstract and/or introduction should clearly state the claims made, including the contributions  
589 made in the paper and important assumptions and limitations. A No or NA answer to this  
590 question will not be perceived well by the reviewers.
- 591 • The claims made should match theoretical and experimental results, and reflect how much the  
592 results can be expected to generalize to other settings.
- 593 • It is fine to include aspirational goals as motivation as long as it is clear that these goals are not  
594 attained by the paper.

### 595 **2. Limitations**

596 Question: Does the paper discuss the limitations of the work performed by the authors?

597 Answer: [\[Yes\]](#)

598 Justification: The paper acknowledges the dependency on reasonable domain knowledge and coverage  
599 please refer to sec 6

600 Guidelines:

- 601 • The answer NA means that the paper has no limitation while the answer No means that the paper  
602 has limitations, but those are not discussed in the paper.
- 603 • The authors are encouraged to create a separate "Limitations" section in their paper.
- 604 • The paper should point out any strong assumptions and how robust the results are to violations of  
605 these assumptions (e.g., independence assumptions, noiseless settings, model well-specification,  
606 asymptotic approximations only holding locally). The authors should reflect on how these  
607 assumptions might be violated in practice and what the implications would be.
- 608 • The authors should reflect on the scope of the claims made, e.g., if the approach was only tested  
609 on a few datasets or with a few runs. In general, empirical results often depend on implicit  
610 assumptions, which should be articulated.
- 611 • The authors should reflect on the factors that influence the performance of the approach. For  
612 example, a facial recognition algorithm may perform poorly when image resolution is low or  
613 images are taken in low lighting. Or a speech-to-text system might not be used reliably to provide  
614 closed captions for online lectures because it fails to handle technical jargon.
- 615 • The authors should discuss the computational efficiency of the proposed algorithms and how  
616 they scale with dataset size.
- 617 • If applicable, the authors should discuss possible limitations of their approach to address problems  
618 of privacy and fairness.
- 619 • While the authors might fear that complete honesty about limitations might be used by reviewers  
620 as grounds for rejection, a worse outcome might be that reviewers discover limitations that  
621 aren't acknowledged in the paper. The authors should use their best judgment and recognize  
622 that individual actions in favor of transparency play an important role in developing norms that  
623 preserve the integrity of the community. Reviewers will be specifically instructed to not penalize  
624 honesty concerning limitations.

### 625 **3. Theory Assumptions and Proofs**

626 Question: For each theoretical result, does the paper provide the full set of assumptions and a complete  
627 (and correct) proof?

628 Answer: [\[Yes\]](#)

629 Justification: Please refer to App A in the supplement material for theoretical analysis and proofs.

630 Guidelines:

- 631 • The answer NA means that the paper does not include theoretical results.
- 632 • All the theorems, formulas, and proofs in the paper should be numbered and cross-referenced.



- 633 • All assumptions should be clearly stated or referenced in the statement of any theorems.
- 634 • The proofs can either appear in the main paper or the supplemental material, but if they appear in
- 635 the supplemental material, the authors are encouraged to provide a short proof sketch to provide
- 636 intuition.
- 637 • Inversely, any informal proof provided in the core of the paper should be complemented by
- 638 formal proofs provided in appendix or supplemental material.
- 639 • Theorems and Lemmas that the proof relies upon should be properly referenced.

#### 640 4. Experimental Result Reproducibility

641 Question: Does the paper fully disclose all the information needed to reproduce the main experimental  
642 results of the paper to the extent that it affects the main claims and/or conclusions of the paper  
643 (regardless of whether the code and data are provided or not)?

644 Answer: [Yes]

645 Justification: Yes all hyper-parameters and experimental setting have been clearly listed in the paper.  
646 Please refer to App F and sec 5.1.

647 Guidelines:

- 648 • The answer NA means that the paper does not include experiments.
- 649 • If the paper includes experiments, a No answer to this question will not be perceived well by the
- 650 reviewers: Making the paper reproducible is important, regardless of whether the code and data
- 651 are provided or not.
- 652 • If the contribution is a dataset and/or model, the authors should describe the steps taken to make
- 653 their results reproducible or verifiable.
- 654 • Depending on the contribution, reproducibility can be accomplished in various ways. For
- 655 example, if the contribution is a novel architecture, describing the architecture fully might suffice,
- 656 or if the contribution is a specific model and empirical evaluation, it may be necessary to either
- 657 make it possible for others to replicate the model with the same dataset, or provide access to
- 658 the model. In general, releasing code and data is often one good way to accomplish this, but
- 659 reproducibility can also be provided via detailed instructions for how to replicate the results,
- 660 access to a hosted model (e.g., in the case of a large language model), releasing of a model
- 661 checkpoint, or other means that are appropriate to the research performed.
- 662 • While NeurIPS does not require releasing code, the conference does require all submissions
- 663 to provide some reasonable avenue for reproducibility, which may depend on the nature of the
- 664 contribution. For example
- 665 (a) If the contribution is primarily a new algorithm, the paper should make it clear how to
- 666 reproduce that algorithm.
- 667 (b) If the contribution is primarily a new model architecture, the paper should describe the
- 668 architecture clearly and fully.
- 669 (c) If the contribution is a new model (e.g., a large language model), then there should either be
- 670 a way to access this model for reproducing the results or a way to reproduce the model (e.g.,
- 671 with an open-source dataset or instructions for how to construct the dataset).
- 672 (d) We recognize that reproducibility may be tricky in some cases, in which case authors are
- 673 welcome to describe the particular way they provide for reproducibility. In the case of
- 674 closed-source models, it may be that access to the model is limited in some way (e.g.,
- 675 to registered users), but it should be possible for other researchers to have some path to
- 676 reproducing or verifying the results.

#### 677 5. Open access to data and code

678 Question: Does the paper provide open access to the data and code, with sufficient instructions to  
679 faithfully reproduce the main experimental results, as described in supplemental material?

680 Answer: [Yes]

681 Justification: The code is provided with the submission in a zip file with Readme for instructions.

682 Guidelines:

- 683 • The answer NA means that paper does not include experiments requiring code.
- 684 • Please see the NeurIPS code and data submission guidelines ([https://nips.cc/public/](https://nips.cc/public/guides/CodeSubmissionPolicy)
- 685 [guides/CodeSubmissionPolicy](https://nips.cc/public/guides/CodeSubmissionPolicy)) for more details.
- 686 • While we encourage the release of code and data, we understand that this might not be possible,
- 687 so “No” is an acceptable answer. Papers cannot be rejected simply for not including code, unless
- 688 this is central to the contribution (e.g., for a new open-source benchmark).
- 689 • The instructions should contain the exact command and environment needed to run to reproduce
- 690 the results. See the NeurIPS code and data submission guidelines ([https://nips.cc/public/](https://nips.cc/public/guides/CodeSubmissionPolicy)
- 691 [guides/CodeSubmissionPolicy](https://nips.cc/public/guides/CodeSubmissionPolicy)) for more details.

- 692
- 693
- 694
- 695
- 696
- 697
- 698
- 699
- 700
- The authors should provide instructions on data access and preparation, including how to access the raw data, preprocessed data, intermediate data, and generated data, etc.
  - The authors should provide scripts to reproduce all experimental results for the new proposed method and baselines. If only a subset of experiments are reproducible, they should state which ones are omitted from the script and why.
  - At submission time, to preserve anonymity, the authors should release anonymized versions (if applicable).
  - Providing as much information as possible in supplemental material (appended to the paper) is recommended, but including URLs to data and code is permitted.

## 701 6. Experimental Setting/Details

702 Question: Does the paper specify all the training and test details (e.g., data splits, hyperparameters,  
703 how they were chosen, type of optimizer, etc.) necessary to understand the results?

704 Answer: [Yes]

705 Justification: The paper uses open source code to create the dataset and lists the modifications in  
706 details in the main paper and supplement material. Please refer to App F and sec 5.1.

707 Guidelines:

- 708
- 709
- 710
- 711
- The answer NA means that the paper does not include experiments.
  - The experimental setting should be presented in the core of the paper to a level of detail that is necessary to appreciate the results and make sense of them.
  - The full details can be provided either with the code, in appendix, or as supplemental material.

## 712 7. Experiment Statistical Significance

713 Question: Does the paper report error bars suitably and correctly defined or other appropriate informa-  
714 tion about the statistical significance of the experiments?

715 Answer: [Yes]

716 Justification: All experiments have been run on 3 random seeds and the error bounds have been  
717 reported in Table 5.1, Table 2 and Table 3.

718 Guidelines:

- 719
- 720
- 721
- 722
- 723
- 724
- 725
- 726
- 727
- 728
- 729
- 730
- 731
- 732
- 733
- 734
- 735
- 736
- 737
- The answer NA means that the paper does not include experiments.
  - The authors should answer "Yes" if the results are accompanied by error bars, confidence intervals, or statistical significance tests, at least for the experiments that support the main claims of the paper.
  - The factors of variability that the error bars are capturing should be clearly stated (for example, train/test split, initialization, random drawing of some parameter, or overall run with given experimental conditions).
  - The method for calculating the error bars should be explained (closed form formula, call to a library function, bootstrap, etc.)
  - The assumptions made should be given (e.g., Normally distributed errors).
  - It should be clear whether the error bar is the standard deviation or the standard error of the mean.
  - It is OK to report 1-sigma error bars, but one should state it. The authors should preferably report a 2-sigma error bar than state that they have a 96% CI, if the hypothesis of Normality of errors is not verified.
  - For asymmetric distributions, the authors should be careful not to show in tables or figures symmetric error bars that would yield results that are out of range (e.g. negative error rates).
  - If error bars are reported in tables or plots, The authors should explain in the text how they were calculated and reference the corresponding figures or tables in the text.

## 738 8. Experiments Compute Resources

739 Question: For each experiment, does the paper provide sufficient information on the computer  
740 resources (type of compute workers, memory, time of execution) needed to reproduce the experiments?

741 Answer: [Yes]

742 Justification: Please refer to the Experimental setup section in the main paper sec 5.1.

743 Guidelines:

- 744
- 745
- 746
- The answer NA means that the paper does not include experiments.
  - The paper should indicate the type of compute workers CPU or GPU, internal cluster, or cloud provider, including relevant memory and storage.

- 747
- The paper should provide the amount of compute required for each of the individual experimental runs as well as estimate the total compute.
- 748
- The paper should disclose whether the full research project required more compute than the experiments reported in the paper (e.g., preliminary or failed experiments that didn't make it into the paper).
- 749
- 750
- 751

## 752 9. Code Of Ethics

753 Question: Does the research conducted in the paper conform, in every respect, with the NeurIPS Code  
754 of Ethics <https://neurips.cc/public/EthicsGuidelines?>

755 Answer: [Yes]

756 Justification: The authors have reviewed the code of ethics and the paper adheres to it.

757 Guidelines:

- The answer NA means that the authors have not reviewed the NeurIPS Code of Ethics.
  - If the authors answer No, they should explain the special circumstances that require a deviation from the Code of Ethics.
  - The authors should make sure to preserve anonymity (e.g., if there is a special consideration due to laws or regulations in their jurisdiction).
- 758
- 759
- 760
- 761
- 762

## 763 10. Broader Impacts

764 Question: Does the paper discuss both potential positive societal impacts and negative societal impacts  
765 of the work performed?

766 Answer: [Yes]

767 Justification: Please refer to the section broader impacts 7.

768 Guidelines:

- The answer NA means that there is no societal impact of the work performed.
  - If the authors answer NA or No, they should explain why their work has no societal impact or why the paper does not address societal impact.
  - Examples of negative societal impacts include potential malicious or unintended uses (e.g., disinformation, generating fake profiles, surveillance), fairness considerations (e.g., deployment of technologies that could make decisions that unfairly impact specific groups), privacy considerations, and security considerations.
  - The conference expects that many papers will be foundational research and not tied to particular applications, let alone deployments. However, if there is a direct path to any negative applications, the authors should point it out. For example, it is legitimate to point out that an improvement in the quality of generative models could be used to generate deepfakes for disinformation. On the other hand, it is not needed to point out that a generic algorithm for optimizing neural networks could enable people to train models that generate Deepfakes faster.
  - The authors should consider possible harms that could arise when the technology is being used as intended and functioning correctly, harms that could arise when the technology is being used as intended but gives incorrect results, and harms following from (intentional or unintentional) misuse of the technology.
  - If there are negative societal impacts, the authors could also discuss possible mitigation strategies (e.g., gated release of models, providing defenses in addition to attacks, mechanisms for monitoring misuse, mechanisms to monitor how a system learns from feedback over time, improving the efficiency and accessibility of ML).
- 769
- 770
- 771
- 772
- 773
- 774
- 775
- 776
- 777
- 778
- 779
- 780
- 781
- 782
- 783
- 784
- 785
- 786
- 787
- 788
- 789

## 790 11. Safeguards

791 Question: Does the paper describe safeguards that have been put in place for responsible release of  
792 data or models that have a high risk for misuse (e.g., pretrained language models, image generators, or  
793 scraped datasets)?

794 Answer: [NA]

795 Justification: The algorithm proposed in this paper does not pose any such risk of misuse.

796 Guidelines:

- The answer NA means that the paper poses no such risks.
  - Released models that have a high risk for misuse or dual-use should be released with necessary safeguards to allow for controlled use of the model, for example by requiring that users adhere to usage guidelines or restrictions to access the model or implementing safety filters.
  - Datasets that have been scraped from the Internet could pose safety risks. The authors should describe how they avoided releasing unsafe images.
- 797
- 798
- 799
- 800
- 801
- 802

803 • We recognize that providing effective safeguards is challenging, and many papers do not require  
804 this, but we encourage authors to take this into account and make a best faith effort.

## 805 12. Licenses for existing assets

806 Question: Are the creators or original owners of assets (e.g., code, data, models), used in the paper,  
807 properly credited and are the license and terms of use explicitly mentioned and properly respected?

808 Answer: [Yes]

809 Justification: All codes and datasets used in this paper are under MIT licence and the original owners  
810 have been cited.

811 Guidelines:

- 812 • The answer NA means that the paper does not use existing assets.
- 813 • The authors should cite the original paper that produced the code package or dataset.
- 814 • The authors should state which version of the asset is used and, if possible, include a URL.
- 815 • The name of the license (e.g., CC-BY 4.0) should be included for each asset.
- 816 • For scraped data from a particular source (e.g., website), the copyright and terms of service of  
817 that source should be provided.
- 818 • If assets are released, the license, copyright information, and terms of use in the package should  
819 be provided. For popular datasets, [paperswithcode.com/datasets](https://paperswithcode.com/datasets) has curated licenses for  
820 some datasets. Their licensing guide can help determine the license of a dataset.
- 821 • For existing datasets that are re-packaged, both the original license and the license of the derived  
822 asset (if it has changed) should be provided.
- 823 • If this information is not available online, the authors are encouraged to reach out to the asset's  
824 creators.

## 825 13. New Assets

826 Question: Are new assets introduced in the paper well documented and is the documentation provided  
827 alongside the assets?

828 Answer: [NA]

829 Justification: No new assets have been introduced in this paper.

830 Guidelines:

- 831 • The answer NA means that the paper does not release new assets.
- 832 • Researchers should communicate the details of the dataset/code/model as part of their sub-  
833 missions via structured templates. This includes details about training, license, limitations,  
834 etc.
- 835 • The paper should discuss whether and how consent was obtained from people whose asset is  
836 used.
- 837 • At submission time, remember to anonymize your assets (if applicable). You can either create an  
838 anonymized URL or include an anonymized zip file.

## 839 14. Crowdsourcing and Research with Human Subjects

840 Question: For crowdsourcing experiments and research with human subjects, does the paper include  
841 the full text of instructions given to participants and screenshots, if applicable, as well as details about  
842 compensation (if any)?

843 Answer: [NA]

844 Justification: The paper did not require any crowdsourcing or human subject for experimentation.

845 Guidelines:

- 846 • The answer NA means that the paper does not involve crowdsourcing nor research with human  
847 subjects.
- 848 • Including this information in the supplemental material is fine, but if the main contribution of the  
849 paper involves human subjects, then as much detail as possible should be included in the main  
850 paper.
- 851 • According to the NeurIPS Code of Ethics, workers involved in data collection, curation, or other  
852 labor should be paid at least the minimum wage in the country of the data collector.

## 853 15. Institutional Review Board (IRB) Approvals or Equivalent for Research with Human Subjects

854 Question: Does the paper describe potential risks incurred by study participants, whether such  
855 risks were disclosed to the subjects, and whether Institutional Review Board (IRB) approvals (or an  
856 equivalent approval/review based on the requirements of your country or institution) were obtained?

857

Answer: [NA]

858

Justification: The paper did not require any crowdsourcing or human subject for experimentation.

859

Guidelines:

860

- The answer NA means that the paper does not involve crowdsourcing nor research with human subjects.

861

862

- Depending on the country in which research is conducted, IRB approval (or equivalent) may be required for any human subjects research. If you obtained IRB approval, you should clearly state this in the paper.

863

864

865

- We recognize that the procedures for this may vary significantly between institutions and locations, and we expect authors to adhere to the NeurIPS Code of Ethics and the guidelines for their institution.

866

867

868

- For initial submissions, do not include any information that would break anonymity (if applicable), such as the institution conducting the review.

869