
Simple Data Sharing for Multi-Tasked Goal-Oriented Problems

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Many important sequential decision problems – from robotics, games to logistics
2 – are multi-tasked and goal-oriented. In this work, we frame them as Contextual
3 Goal Oriented (CGO) problems, a goal-reaching special case of the contextual
4 Markov decision process. CGO is a framework for designing multi-task agents
5 that can follow instructions (represented by contexts) to solve goal-oriented tasks.
6 We show that CGO problem can be systematically tackled using datasets that are
7 commonly obtainable: an unsupervised interaction dataset of transitions and a su-
8 pervised dataset of context-goal pairs. Leveraging the goal-oriented structure of
9 CGO, we propose a simple data sharing technique that can provably solve CGO
10 problems offline under natural assumptions on the datasets’ quality. While an of-
11 fline CGO problem is a special case of offline reinforcement learning (RL) with
12 unlabelled data, running a generic offline RL algorithm here can be overly con-
13 servative since the goal-oriented structure of CGO is ignored. In contrast, our
14 approach carefully constructs an augmented Markov Decision Process (MDP) to
15 avoid introducing unnecessary pessimistic bias. In the experiments, we demon-
16 strate our algorithm can learn near-optimal context-conditioned policies in simu-
17 lated CGO problems, outperforming offline RL baselines.

18 1 Introduction

19 Goal-Oriented (GO) problems (Kaelbling, 1993) are an important class of sequential decision-
20 making problems with widespread applications, ranging from robotics (Yu & Mooney, 2023) to
21 game-playing (Hessel et al., 2019) to real-world logistics (Mirowski et al., 2018). Many of these
22 problems are multi-tasked: rather than aiming toward a single goal, the agent needs to reach task-
23 specific goals based on the task instruction it receives. In this work, we frame these multi-tasked
24 goal-oriented applications as Contextual GO (CGO) problems and design a simple algorithm that
25 can provably solve them using offline datasets that are commonly available in CGO applications.

26 CGO problem is a special case of contextual Markov Decision Process (MDP) (Hallak et al., 2015).
27 In a CGO problem, each task is a reaching problem with a goal set that is communicated indirectly
28 to the agent via a context. CGO problem includes the classical GO problem as a special case,
29 where the context is just the target goal, but in general contexts in CGO problem can convey rich,
30 high-level task instructions. In robotics, e.g., common contexts are verbal instructions like “clean
31 up the table” whereas goals are specific configurations (e.g., a clean table) in the environment. In
32 games, contexts can be side-quests for the player to accomplish, and in logistics contexts describe
33 origins and destinations of journeys an operator should execute. We will use navigation as a running
34 example in this paper. Imagine instructing a truck operator with the context “Deliver goods to a
35 warehouse in the Bay area”. Given the context, they must first infer a goal (e.g., a warehouse
36 location) and implement a policy to efficiently navigate to the goal.

37 CGO problems are challenging, because the rewards are sparse (non-zero rewards only when reach-
38 ing goals) and the contexts can be difficult to interpret into feasible goals. However, CGO problem
39 has an important structure that the transition dynamics (e.g., navigating a city road network) are
40 independent of the context (e.g., journey origin and destination), and efficient multitask learning can
41 be achieved by sharing dynamics data across tasks or contexts.

42 We study offline Reinforcement Learning (RL) for CGO problems. Offline learning is timely for
43 CGO problems given the recent availability of suitable massive datasets. We identify two different
44 kinds of datasets that are commonly available in CGO applications – an (unsupervised) *dynamics*
45 dataset of agent trajectories, and a (supervised) *context-goal* dataset of pairs of contexts and goals.
46 In robotics, task-agnostic play data can be obtained at scale (Lynch et al., 2020; Walke et al., 2023)
47 in an unsupervised manner whereas instruction datasets (e.g., Misra et al. (2016)) allow supervised
48 learning of the context-goal mapping. In navigation, self-driving car trajectories (e.g., Wilson et al.
49 (2021); Sun et al. (2020)) allow us to learn dynamics whereas landmarks datasets (e.g. Mirowski
50 et al. (2018); Hahn et al. (2021)) allow us to map the contexts to goals.

51 We propose a Simple Data Sharing (SDS) technique that can provably solve CGO problems subject
52 to natural assumptions on the datasets’ quality. We prove that SDS can learn a near-optimal policy
53 for the CGO problem with high probability, as long as the distribution generating the context-goal
54 dataset covers the target context and the distribution generating the dynamics dataset covers a feasi-
55 ble path to the target goal set. SDS is a reduction-based technique that can be implemented on top of
56 a standard offline RL algorithm. Our key insight is to carefully construct an action-augmented MDP
57 such that the dynamics dataset and context-goal dataset can be reconciled together as a standard
58 reward-labeled offline dataset.

59 To our knowledge, SDS is the first offline algorithm that can provably solve CGO problems with
60 just positive data (i.e., the context-goal dataset). While the offline CGO problem here can be cast as
61 an offline RL problem with unlabeled data (i.e., viewing each {context, state} pair as a composite
62 state¹), existing theoretical results (Yu et al., 2022; Hu et al., 2023; Li et al., 2023a) indicate that both
63 positive data and negative data (i.e., pairs of context and non-goal data) are needed.² An alternative
64 approach to offline CGO problems is to predict goals based on contexts and then run offline goal-
65 conditioned RL (Ma et al., 2022). This approach only needs positive data in learning the predictor,
66 but it can fail when the predicted goal is not reachable from the initial state. In the truck operator
67 example, suppose that there are two warehouses on either side of a river but the bridge across the
68 river is closed to traffic. The goal predictor must reason about the connectivity of the road network
69 when it sets goals; otherwise it may set an infeasible goal (e.g., a warehouse on the other side of the
70 river) that no goal-conditioned policy can successfully execute.

71 We contribute an effective SDS technique and a new analysis technique that formally proves that
72 CGO problem can be solved offline with just dynamics data and context-goal data (i.e. positive
73 data), without the need of negative data. We also show that SDS can be implemented on top of
74 existing offline RL algorithms (with concrete instantiations for PSPI (Xie et al., 2021) in Section 3.3
75 and IQL (Kostrikov et al., 2021) in Section 4). In addition to theoretical analyses, we conduct several
76 experiments in simulated domains, confirming that SDS outperforms SOTA offline RL baselines
77 designed for unlabeled data. Finally, we situate our contributions within the vast literature on Goal-
78 Oriented RL (Kaelbling, 1993) and contextual MDPs (Hallak et al., 2015) in Appendix A.

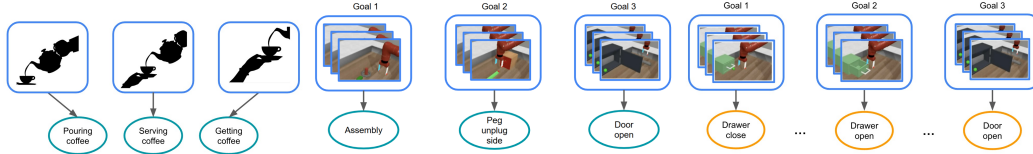
79 2 Preliminaries

80 2.1 Contextual Goal-Oriented (CGO) Problem

81 A Contextual Goal-Oriented (CGO) problem describes a multi-task goal-oriented setting with a
82 *shared* transition kernel. We consider a Markovian CGO problem with an infinite horizon, defined
83 by the tuple $\mathcal{M} = (\mathcal{S}, \mathcal{A}, P, R, \gamma, \mathcal{C}, d_0)$, where \mathcal{S} is the state space, \mathcal{A} is the action space, $P : \mathcal{S} \times \mathcal{A} \rightarrow \Delta(\mathcal{S})$ is the transition kernel, $R : \mathcal{S} \times \mathcal{C} \rightarrow \{0, 1\}$ is the reward function, $\gamma \in [0, 1)$
84 is the discount factor, \mathcal{C} is the context space, and finally Δ denotes the space of distributions. We
85

¹Context-goal data can be processed into reward-labeled data, whereas dynamics data from the original MDP imputed with all of the contexts seen in the context-goal dataset becomes the reward-unlabeled data.

²Additionally reward-labeled data covering the full trajectory is necessary for general offline RL. But for GO problems, we show that a weaker condition of covering only the goals is sufficient. Existing algorithms for offline RL with unlabeled data may work with this weaker notion of coverage, but it is unclear how to prove it.



(a) Similar goal sets with different contexts (b) Distinct goal sets with different but small number of contexts (c) Overlapping goal sets across contexts but with an empty intersection

Figure 1: The interplay between contexts and goals in a Contextual Goal-Oriented (CGO) problem characterizes many real-world multi-task settings. (a) All the contexts may share similar goal sets (e.g., pouring coffee). (b) Each context may map to different goal sets (e.g., general-purpose robotics). (c) Contexts may have different overlapping goal sets, creating a complex CGO problem.

do not assume any particular topology on \mathcal{S} , \mathcal{A} and \mathcal{C} and they can be continuous. Each context $c \in \mathcal{C}$ specifies a goal-reaching task with a goal set $G_c \subset \mathcal{S}$, and reaching any goal in the goal set G_c is regarded as successful. The reward function is hence defined as $R(s, c) = \mathbb{1}(s \in G_c)$. An episode of a CGO problem starts from an initial state s_0 and a context c sampled according to a distribution $d_0(s_0, c)$, and it terminates when the agent reaches the goal set G_c . During the episode, c does not change; only s_t changes (according to $P(s'|s, a)$) and the transition kernel $P(s'|s, a)$ is context independent. The classical GO problem (Kaelbling, 1993) is a special case of CGO, where a multi-goal problem can be viewed as multiple contexts with each context describing a goal.

Spectrum of CGO Problem Figure 1 illustrates different CGO problems encountered when learning a language-conditioned control policy for a robot manipulator. s describes the robot and the world state, a is the robot action, and c is the language instruction. For each instruction c , the manipulation task for the robot is a reaching problem to a set of targeted robot and world states. The simplest CGO instance is when most of the contexts $c \in \mathcal{C}$ correspond to the very similar goal sets, as shown in Figure 1a. In this case, a context-agnostic policy can be near-optimal³. When different contexts have non-overlapping goal sets G_c and the number of contexts are small (as in Figure 1b), the problem is essentially multi-task RL which *requires* context-conditioned policies. In its full complexity, the number of contexts can be infinite; and goal sets of different contexts could overlaps while their intersection is empty, as shown in Figure 1c. A CGO agent thus needs to learn how to respond to different contexts as well as transfer knowledge efficiently across contexts.

Objective Since the context carries rich information, a CGO policy in general is context-conditioned, i.e., $\pi : \mathcal{S} \times \mathcal{C} \rightarrow \Delta(\mathcal{A})$. The performance of a policy π is measured by its return, $J(\pi) := \mathbb{E}_{\pi, P, d_0} \left[\sum_{t=0}^T \gamma^t R(s_t, c) \right]$, where T is the time the agent first enters G_c (a random variable dependent on π , P and d_0), and \mathbb{E}_{π, P, d_0} denotes the expectation over trajectories generated by running π with P starting from s_0 , c sampled from d_0 . We can view the return as the average success rate of reaching *any* goal in the goal set G_c when the problem horizon is exponentially distributed (according to the discount γ). A CGO algorithm takes a policy class $\Pi = \{\pi : \mathcal{S} \times \mathcal{C} \rightarrow \Delta(\mathcal{A})\}$ as input and returns a near-optimal policy π^\dagger such that $J(\pi^\dagger) \approx \max_{\pi \in \Pi} J(\pi)$.

2.2 Offline Learning

We aim to solve CGO problems using offline datasets without additional online environment interactions, à la offline RL. We identify two types of data that are commonly available: $D_{\text{dyn}} := \{(s, a, s')\}$ is an *unsupervised* dataset of agent trajectories collected from $P(s'|s, a)$, whereas $D_{\text{goal}} := \{(c, s) : s \in G_c\}$ is a *supervised* dataset of context-goal pairs. Different offline CGO algorithms can be judged based on the assumptions they require on $\{D_{\text{dyn}}, D_{\text{goal}}\}$, such as what the datasets should cover and how much data are needed to learn π^\dagger . No algorithm, to our knowledge, can *provably* learn near-optimal π^\dagger using *only* the positive D_{goal} data (i.e., without needing additional *negative* data of non-goal examples) when combined with D_{dyn} data. In the next section, we demonstrate how to leverage the special structure of the CGO problem to design provably correct offline algorithms. This insight leads to a Simple Data Sharing (SDS) scheme that can enable existing offline

³Indeed we show in Section 4 that some existing multi-task RL benchmarks are in this regime where a context-agnostic Implicit Q-Learning (IQL) (Kostrikov et al., 2021) baseline performs well.

124 RL algorithms (designed for fully labeled data) to solve offline CGO problems using *just* the positive
 125 goal-labeled data without needing any additional non-goal examples, or reward learning.

126 2.3 Notation and Assumption

127 Before presenting the main results, we introduce some definitions and shorthand to make the pre-
 128 sentation more readable. We introduce a fictitious zero-reward absorbing state s^+ and modify the
 129 dynamics such that whenever the agent enters G_c it transits to s^+ in the next time step (for all ac-
 130 tions) and stays there forever. This is a standard technique to convert a goal reaching problem (with
 131 a random problem horizon) to an infinite horizon problem. It does *not* change the problem.

132 Specifically, we extend the reward and the dynamics as follows: We define $\bar{\mathcal{S}} = \mathcal{S} \cup \{s^+\}$, $\mathcal{X} :=$
 133 $\bar{\mathcal{S}} \times \mathcal{C}$, and $\bar{\mathcal{X}} := \bar{\mathcal{S}} \times \mathcal{C}$. In addition, we define $\mathcal{X}^+ := \{x : x = (s, c), s = s^+, c \in \mathcal{C}\}$. We use
 134 G to denote the goal set on \mathcal{X} , i.e., $G := \{x \in \mathcal{X} : x = (s, c), s \in G_c\}$. With abuse of notation,
 135 we define the reward function and the transition kernel on \mathcal{X} accordingly as $R(x) = \mathbb{1}(s \in G_c)$
 136 and $P(x'|x, a) := P(s'|s, c, a)\mathbb{1}(c' = c)$, where $P(s'|s, c, a) := \mathbb{1}(s' = s^+)$ if $s \in G_c$ or $s = s^+$;
 137 otherwise $P(s'|s, c, a) = P(s'|s, a)$, where $x = (s, c)$ and $x' = (s', c')$. Notice the context does not
 138 change in the transition. For all value functions, we define their value at s^+ as zero.

139 Given a policy $\pi : \mathcal{X} \rightarrow \Delta(\mathcal{A})$, we define its state-action value function (i.e., Q function) as
 140 $Q^\pi(x, a) := \mathbb{E}_{\pi, P} [\sum_{t=0}^{\infty} \gamma^t R(x) | x_0 = x, a_0 = a]$. We use $V^\pi(x) := Q^\pi(x, \pi)$ to denote the value
 141 function π , where $f(\pi) := \mathbb{E}_{a \sim \pi} [f(a)]$. By construction, we have $Q^\pi(x, a), V^\pi(x) \in [0, 1], \forall x \in$
 142 $\mathcal{X}, a \in \mathcal{A}$. By these definitions, we can write the return $J(\pi) = V^\pi(d_0) = Q^\pi(d_0, \pi)$. We denote
 143 π^* as the optimal policy and define $Q^* := Q^{\pi^*}, V^* := V^{\pi^*}$.

144 **Data Assumption** We suppose that there are two distributions $\mu_{\text{dyn}}(s, a, s')$ and $\mu_{\text{goal}}(s, c)$, where
 145 $\mu_{\text{dyn}}(s'|s, a) = P(s'|s, a)$ and μ_{goal} has support within G_c , i.e., $\mu_{\text{goal}}(s|c) > 0 \Rightarrow s \in G_c$. We
 146 assume that D_{dyn} and D_{goal} are i.i.d. samples drawn from μ_{dyn} and μ_{goal} , i.e.,

$$D_{\text{dyn}} = \{(s_i, a_i, s'_i) \sim \mu_{\text{dyn}}\} \quad \text{and} \quad D_{\text{goal}} = \{(s_j, c_j) \sim \mu_{\text{goal}}\}.$$

147 We suppose that $x \sim d_0$ is not in G almost surely. This is to simplify the presentation. If $x \in G$, the
 148 agent reaches its goal immediately and no learning is needed.

149 3 Simple Data Sharing To Solve CGO Problems

150 The key idea of SDS is the construction of an *action*-augmented MDP with which the dynamics
 151 and context-goal datasets can be combined into a conventional offline RL dataset. In the following,
 152 first we describe this action-augmented MDP (Section 3.1) and show that it preserves the optimal
 153 policies of the original MDP (Appendix B.1). We then outline a practical algorithm to convert the
 154 two datasets of an offline CGO problem into a dataset for this augmented MDP (Section 3.2) such
 155 that any generic offline RL algorithm can be used as a solver. Finally, in Section 3.3, we theoretically
 156 analyze an instantiation of SDS based on PSPI (Xie et al., 2021) and show that SDS can provably
 157 find a near-optimal policy for the CGO problem.

158 3.1 Action-Augmented MDP

159 One reason why offline RL cannot directly leverage D_{dyn} and D_{goal} to solve a CGO problem is that
 160 each goal-reaching problem has its own context-specific termination criterion. Notice that although
 161 the dynamics datasets D_{dyn} is consistent with the original MDP transition kernel (i.e. $P(s'|s, a)$),
 162 it is however not consistent with the transition kernel $P(x'|x, a)$ (which also includes the effect of
 163 context-specific termination) of the context-augmented MDP in Section 2.3. This is easiest to see
 164 if some $s \in G_c$ in the D_{goal} dataset is also observed in the dynamics dataset. D_{dyn} will imply from
 165 (s, a, s') that action a can transition to s' , however D_{goal} implies that all actions at s will transition to
 166 s^+ . This conflict means that combining the two datasets naively leads to an inconsistent algorithm.

167 We propose a new augmented MDP, which augments the action space of the context-augmented
 168 MDP in Section 2.3 with a fictitious action a^+ to avoid conflicts across D_{dyn} and D_{goal} . Define
 169 $\bar{\mathcal{A}} = \mathcal{A} \cup \{a^+\}$. The reward in this action-augmented MDP is now *action-dependent*, for $x =$
 170 $(s, c) \in \mathcal{X}$, $\bar{R}(x, a) := \mathbb{1}(s \in G_c)\mathbb{1}(a = a^+)$ and the transition upon taking action a^+ is defined as
 171 $P(x'|x, a^+) := \mathbb{1}(s' = s^+)$ and $P(x'|x, a) := P(s'|s, a)\mathbb{1}(c' = c)$ for other actions.

Algorithm 1 Simple Data Sharing (SDS) for CGO

Input: Dynamics dataset D_{dyn} , context-goal dataset D_{goal}

for each sample $(s, c) \sim D_{\text{goal}}$ **do**

 Create transition⁴ $(x, a^+, 1, x^+)$, where $x = (s, c)$ and $x^+ = (s^+, c)$, add it to \bar{D}_{goal}

end for

for each $(s, a, s') \sim D_{\text{dyn}}$ **do**

for each $(\cdot, c) \sim D_{\text{goal}}$ **do**

 Create transition $(x, a^+, 0, x')$, where $x = (s, c)$ and $x' = (s', c)$, add it to \bar{D}_{dyn}

end for

end for

Output: \bar{D}_{dyn} and \bar{D}_{goal}

172 We denote this action-augmented MDP as $\bar{\mathcal{M}} := (\bar{\mathcal{X}}, \bar{\mathcal{A}}, \bar{R}, \bar{P}, \gamma)$. For policy $\pi : \mathcal{X} \rightarrow \Delta(\mathcal{A})$ and
173 value functions $f : \mathcal{X} \times \mathcal{A} \rightarrow [0, 1]$ defined in the original MDP, we define their extensions on $\bar{\mathcal{M}}$:

$$\bar{\pi}(a|x) = \begin{cases} \pi(a|x), & x \notin G \\ a^+, & \text{otherwise} \end{cases} \quad \text{and} \quad \bar{f}_g(x, a) = \begin{cases} g(x), & a = a^+ \text{ and } x \notin \mathcal{X}^+ \\ 0, & x \in \mathcal{X}^+ \\ f(x, a), & \text{otherwise} \end{cases}$$

174 where the extension of f is based on a function $g : \mathcal{X} \rightarrow [0, 1]$ which determines its value at a^+ .

175 We show in Appendix B.1 (see Lemma B.3) that the regret of a policy extended to the augmented
176 MDP is equal to the regret of the policy in the original MDP describing the CGO problem, and
177 any policy defined in the augmented MDP can be converted into that in the original MDP without
178 increasing the regret. Thus, solving the augmented MDP can yield correspondingly optimal policies
179 for the original problem. We next sketch a practical technique to combine D_{dyn} and D_{goal} along with
180 the fictitious action labels a^+ such that we can solve the action-augmented MDP effectively.

181 3.2 Practical Algorithm: Simple Data Sharing

182 In Algorithm 1 we sketch our Simple Data Sharing (SDS) technique. It takes the two datasets
183 D_{dyn} and D_{goal} as input, and produces a single dataset $\bar{D}_{\text{dyn}} \cup \bar{D}_{\text{goal}}$ that is suitable for use by any
184 offline RL algorithm like CQL (Kumar et al., 2020), IQL (Kostrikov et al., 2021), PSPI (Xie et al.,
185 2021), ATAC (Cheng et al., 2022) etc. Notice that any policy returned by the offline RL algorithm
186 can be executed in the CGO problem by simply masking out the a^+ action. We note that in practice
187 Algorithm 1 can be implemented as a pre-processing step in the minibatch sampling of a deep offline
188 RL algorithm (as opposed to computing the full \bar{D}_{dyn} and \bar{D}_{goal} once before learning). Empirically,
189 we found that equally balancing the samples \bar{D}_{dyn} and \bar{D}_{goal} generates the best result. Below we
190 analyze SDS theoretically by applying SDS to PSPI (Xie et al., 2021); later in Section 4, we apply
191 SDS to IQL (Kostrikov et al., 2021) in simulation experiments.

192 3.3 Analysis of SDS+PSPI: Information Theoretic Guarantee

193 In this section, we show a formal analysis for our reduction approach, when instantiated with PSPI
194 (Xie et al., 2021). We summarize the main theoretical result as follows.

195 **Theorem 3.1.** *Let π^\dagger denote the learned policy of SDS + PSPI with datasets D_{dyn} and D_{goal} , using
196 value function classes⁵ $\mathcal{F} = \{\mathcal{X} \times \mathcal{A} \rightarrow [0, 1]\}$ and $\mathcal{G} = \{\mathcal{X} \rightarrow [0, 1]\}$. Under realizability and
197 completeness assumptions below, with probability $1 - \delta$, it holds, for any $\pi \in \Pi$,*

$$J(\pi) - J(\pi^\dagger) \leq \mathfrak{C}_{\text{dyn}}(\pi) \sqrt{\epsilon_{\text{dyn}}} + \mathfrak{C}_{\text{goal}}(\pi) \sqrt{\epsilon_{\text{goal}}}$$

198 where $\epsilon_{\text{dyn}} = O\left(\frac{\log(|\mathcal{F}||\mathcal{G}||\Pi|/\delta)}{|D_{\text{dyn}}|}\right)$ and $\epsilon_{\text{goal}} = O\left(\frac{\log(|\mathcal{G}|/\delta)}{|D_{\text{goal}}|}\right)$ are statistical errors, and $\mathfrak{C}_{\text{dyn}}(\pi)$ and
199 $\mathfrak{C}_{\text{goal}}(\pi)$ are concentrability coefficients which decrease as the data coverage increases.

200 **Assumption 3.2** (Realizability). *We assume for any $\pi \in \Pi$, $Q^\pi \in \mathcal{F}$ and $R \in \mathcal{G}$.*

⁴ s^+ is implemented as `terminal=True`.

⁵We state a more general result for non-finite function classes in the appendix.

201 **Assumption 3.3** (Completeness). *We assume: For any $f \in \mathcal{F}$ and $g \in \mathcal{G}$, $\max(g(x), f(x, \pi)) \in \mathcal{F}$;*
 202 *And for any $f \in \mathcal{F}$, $\pi \in \Pi$, $\mathcal{T}^\pi f(x, a) \in \mathcal{F}$, where \mathcal{T}^π is a zero-reward Bellman backup operator*
 203 *with respect to $P(s'|s, a)$: $\mathcal{T}^\pi f(x, a) := \gamma \mathbb{E}_{x' \sim P(s'|s, a) \mathbb{1}(c'=c)} [f(x', \pi)]$.*

204 **Definition 3.4.** *We define the generalized concentrability coefficients:*

$$\mathfrak{C}_{\text{dyn}}(\pi) := \max_{f, f' \in \mathcal{F}} \frac{\|f - \mathcal{T}^\pi f'\|_{\rho_{\bar{c}}^\pi}^2}{\|f - \mathcal{T}^\pi f'\|_{\mu_{\text{dyn}}}^2} \quad \text{and} \quad \mathfrak{C}_{\text{goal}}(\pi) := \max_{g \in \mathcal{G}} \frac{\|g - r\|_{\rho_{\bar{c}}^\pi}^2}{\|g - r\|_{\mu_{\text{goal}}}^2}$$

205 where $\|h\|_\mu^2 := \mathbb{E}_{x \sim \mu} [h(x)^2]$, $\rho_{\bar{c}}^\pi(x, a) = \mathbb{E}_{\pi, P} \left[\sum_{t=0}^{T-1} \gamma^t \mathbb{1}(x_t = x, a_t = a) \right]$, $\rho_{\bar{c}}^\pi(x) =$
 206 $\mathbb{E}_{\pi, P} [\gamma^T \mathbb{1}(x_T = x)]$, and T is the first time the agent enters the goal set.

207 Concentrability coefficients is a generalization notion of density ratio; it describes how much the
 208 (unnormalized) distribution in the numerator is covered by that in the denominator in terms of the
 209 generalization ability of function approximators (Xie et al., 2021). By setting $\pi = \pi^*$ in Theo-
 210 rem 3.1, we see that the policy learned by SDS+PSPI has a small regret as long as the dynamics data
 211 D_{dyn} covers the trajectory of the optimal policy, and the context-goal dataset D_{goal} covers goals the
 212 optimal policy would reach. In other words, SDS+PSPI can provably learn with only the positive
 213 data (i.e., the context-goal dataset) without the need of additional labeling of non-goal samples.

214 **Remark 3.5.** *MAHALO (Li et al., 2023a) is a SOTA offline RL algorithm that can provably learn*
 215 *from unlabeled data. MAHALO can also be implemented on top of PSPI; however, their theoretical*
 216 *result (Theorem D.1) requires a stronger version concentrability, $\max_{g \in \mathcal{G}} \|g - r\|_{\rho_{\bar{c}}^\pi}^2 / \|g - r\|_{\mu_{\text{goal}}}^2$, to be*
 217 *small. In other words, it needs additional labeling of non-goal states.*

218 3.3.1 Algorithm: SDS+PSPI

219 Here we briefly summarize how SDS+PSPI is implemented, without taking literally a^+ and s^+ in
 220 function approximators. Due to space constraints, we defer the details to Appendix B.

221 We consider the information theoretic version of PSPI (Xie et al., 2021) which can be summarized
 222 as follows: For an MDP $(\mathcal{X}, \mathcal{A}, R, P, \gamma)$, given a tuple dataset $D = \{(x, a, r, x')\}$, a policy class Π ,
 223 and a value class \mathcal{F} , it finds the policy through solving the two-player game:

$$\max_{\pi \in \Pi} \min_{f \in \mathcal{F}} f(d_0, \pi) \quad \text{s.t.} \quad \ell(f, f; \pi, D) - \min_{f' \in \mathcal{F}} \ell(f', f; \pi, D) \leq \epsilon_b \quad (1)$$

224 where $f(d_0, \pi) = \mathbb{E}_{x_0 \sim d_0} [f(x_0, \pi)]$, $\ell(f, f'; \pi, D) := \frac{1}{|D|} \sum_{(x, a, r, x') \in D} (f(x, a) - r - f'(x', \pi))^2$.
 225 The term $\ell(f, f; \pi, D) - \min_{f'} \ell(f', f; \pi, D)$ is an empirical estimation of the Bellman error on f
 226 of π on the data distribution μ , i.e. $\mathbb{E}_{x, a \sim \mu} [(f(x, a) - \mathcal{T}^\pi f(x, a))^2]$. It constrains the Bellman error
 227 to be a small ϵ_b , since $\mathbb{E}_{x, a \sim \mu} [(Q^\pi(x, a) - \mathcal{T}^\pi Q^\pi(x, a))^2] = 0$.

228 **Instantiating PSPI** In order to run PSPI on the augmented MDP, we extend the policy class to $\bar{\Pi}$
 229 and define an extended value class $\bar{\mathcal{F}}_{\mathcal{G}}$ based on \mathcal{F} and \mathcal{G} as discussed in Section 3.1. Then we rewrite
 230 the squared Bellman error on the two data distributions ⁶ using equation 6 and Proposition B.4 as:

$$\mathbb{E}_{x, a \sim \mu_{\text{dyn}}} [(\bar{Q}^{\bar{\pi}}(x, a) - \bar{\mathcal{T}}^{\bar{\pi}} \bar{Q}^{\bar{\pi}}(x, a))^2] = \mathbb{E}_{x, a \sim \mu_{\text{dyn}}} [(\bar{Q}^{\bar{\pi}}(x, a) - \gamma \mathbb{E}_{x' \sim \bar{P}(\cdot|x, a)} [\max(R(x'), Q^\pi(x', \pi))])^2]$$

$$\mathbb{E}_{x, a \sim \mu_{\text{goal}}} [(\bar{Q}^{\bar{\pi}}(x, a) - \bar{\mathcal{T}}^{\bar{\pi}} \bar{Q}^{\bar{\pi}}(x, a))^2] = \mathbb{E}_{x, a \sim \mu_{\text{goal}}} [(\bar{Q}^{\bar{\pi}}(x, a^+) - 1)^2]$$

232 where $\bar{\mathcal{T}}^{\bar{\pi}}$ denotes the Bellman backup operator and $\bar{Q}^{\bar{\pi}}$ denotes the Q-function of $\bar{\pi}$ in $\bar{\mathcal{M}}$.

233 Using this expression for the squared Bellman error, we can reformulate the empirical losses in
 234 equation 1:

$$\ell_{\text{dyn}}(\bar{f}_g, \bar{f}'_g; \bar{\pi}) := \frac{1}{|\bar{D}_{\text{dyn}}|} \sum_{(x, a, r, x') \in \bar{D}_{\text{dyn}}} (f(x, a) - \gamma \max(g'(x'), f'(x', \pi)))^2 \quad (2)$$

$$\ell_{\text{goal}}(\bar{f}_g) := \frac{1}{|\bar{D}_{\text{goal}}|} \sum_{(x, a, r, x') \in \bar{D}_{\text{goal}}} (g(x) - 1)^2 \quad (3)$$

235 Using these losses, we can define the two-player game of PSPI for the action-augmented MDP as:

$$\max_{\bar{\pi} \in \bar{\Pi}} \min_{\bar{f}_g \in \bar{\mathcal{F}}} \bar{f}_g(d_0, \bar{\pi}) \quad \text{s.t.} \quad \ell_{\text{dyn}}(\bar{f}_g, \bar{f}_g; \bar{\pi}) - \min_{\bar{f}'_g \in \bar{\mathcal{F}}} \ell_{\text{dyn}}(\bar{f}'_g, \bar{f}_g; \bar{\pi}) \leq \epsilon_{\text{dyn}}, \quad \ell_{\text{goal}}(\bar{f}_g) \leq 0 \quad (4)$$

236 Notice $\bar{f}_g(d_0, \bar{\pi}) = f(d_0, \pi)$, so this problem can be solved using samples without knowing G .

⁶With abuse of notation, we write $\mu_{\text{dyn}}(x, a, x') = \mu_{\text{dyn}}(s, a, s') \mu_{\text{goal}}(c)$ and $\mu_{\text{goal}}(x, a, x') = \mu_{\text{goal}}(c, s) \mathbb{1}(a = a^+) \mathbb{1}(s' = s^+)$. In Algorithm 1, we have $D_{\text{dyn}} \sim \mu_{\text{dyn}}$ and $\bar{D}_{\text{goal}} \sim \mu_{\text{goal}}$.

237 4 Experiments

238 Through experiments we aim to answer the following questions: 1) Does our method work in sce-
239 narios of different context-goal relationships shown in Figure 1, under the data assumptions in Sec-
240 tion 2.3? 2) Under each setting, is there any empirical benefit from using SDS, compared with offline
241 RL baselines (for unlabeled data) that require pessimistic reward learning?

242 4.1 Environments and datasets

243 **Dynamics dataset.** For all experiments, we use the AntMaze-v2 datasets of D4RL (Fu et al., 2020)
244 as dynamics datasets D_{dyn} ; we remove the reward and terminal information labels.

245 **Context-goal dataset.** We construct three levels of context and goal relationships as shown in Fig-
246 ure 1: 1) Figure 1a where multiple contexts define very similar goal sets (Section 4.3); 2) Figure 1b
247 where the number of contexts is finite and the goal sets of different do not overlap (Section 4.4);
248 3) Figure 1c where the contexts are continuous and randomly sampled, the goal sets can overlap
249 but their intersection is empty (Section 4.5). For each environment, we define a context set and
250 an *oracle function* to tell whether a state is within the goal set; this oracle function is only used in
251 data construction and is not accessible to the algorithms tested here. Then given each context, we
252 select states in the dynamics dataset that satisfy the oracle function to construct the goal examples⁷
253 In Appendix E, we include results of the goal set containing samples not from the dynamics dataset.

254 **Evaluation.** Section 4.3, 4.4 and 4.5 contain results where the training and testing contexts are
255 sampled from the same distribution; in Section 4.5 we also test the algorithms with a different
256 context distribution. For evaluation, we use⁸ the oracle function that defines context-goal sets to
257 provide the reward given a certain context in Section 4.4 and 4.5. The evaluation of each context is
258 done by 100 episodes. We train each algorithm for 5 seeds and report the statistics.

259 4.2 Methods

260 Here we describe the algorithms compared in the experiments. To facilitate a clean comparison of
261 different conceptual approaches to solving offline CGO problems, we use IQL (Kostrikov et al.,
262 2021) as the backbone offline algorithm for all the methods. The same set of hyperparameters in
263 IQL is used in all experiments. In the experiments, we use the $-1/0$ reward notion, which can be
264 shown to be the same as the $0/1$ reward notion in terms of ranking policies under the discounted
265 MDP setting. Please see Appendix C.1 for detailed hyperparameters of all methods.

266 **SDS+IQL (Ours).** We apply SDS in Algorithm 1 with IQL as the offline RL algorithm to solve
267 the augmented MDP defined in Section 3.1. More specifically, we set a^+ to be an extra dimension in
268 the action space but mask out extra dimension for policy output. We can think of IQL as optimizing
269 Eq. (2) via expectile regression given the offline dataset..

270 **Reward prediction (RP).** For naive reward prediction, we first convert the context-goal set to
271 a dataset with reward 0 for all $(c, s) \sim D_{\text{goal}}$, and then learn a reward function with the dataset.
272 For policy training, we randomly sample $(s, a, s') \sim D_{\text{dyn}}$ and $c \sim D_{\text{goal}}$ and label the transition
273 with the learned reward: if reward prediction of (c, s') is larger than some threshold, we label the
274 transition with $r = 0$ and `terminal = True`; otherwise we label the transition with $r = -1$ and
275 `terminal = False`. Then we apply IQL with this labeled dataset.

276 **PDS.** For PDS (Hu et al., 2023), we follow the similar procedure as RP but learn a *pessimistic*
277 reward function using ensembles. Then we apply similar steps to label the transitions with contexts
278 and apply IQL with this labeled dataset as RP.

279 **UDS+RP.** On top of RP, we introduce another possible way to learn a reward function while we
280 construct “non-goal” samples in a pessimistic manner: we also sample random $c \sim D_{\text{goal}}$ and
281 $s \sim D_{\text{dyn}}$ and label it with $r = -1$ similar to the spirit of UDS (Yu et al., 2022), then train the
282 reward function with the combined positive and negative dataset. Then we follow the same steps in
283 RP for policy training with the learned reward function.

284 **Context-agnostic IQL.** As discussed in Section 4.1, if we “hack” our context-goal construction
285 method, given contest-goal data we can label the corresponding transitions with `terminal = True`
286 and $r = 0$, and for other transitions and contexts, we label it with `terminal = False` and $r = -1$,

⁷No method in the comparison utilizes this fact.

⁸Exception: in the original AntMaze, we use the D4RL metric, so the results are comparable to the literature.

287 then we will have a labeled offline dataset. We then treat the union of all goal sets as one large goal
 288 set with a single context. It is only to provide a reference to the conventional methods used to solve
 289 AntMaze, but not for comparison with our method or other baselines and **cannot** be implemented in
 290 a real-world offline CGO problem.

291 4.3 Original AntMaze

292 In the original AntMaze, 2D goal locations (contexts) are sampled from a fixed cell in the maze and
 293 perturbed with a small noise, generating very similar goal sets. Our training context set is chosen as
 294 2D locations of the states with terminal=True in the D4RL datasets, and the full state is added as the
 295 goal example. Test contexts and environmental evaluation follow the original AntMaze.

296 **SDS matches the performance of the context-agnostic method under the setting of Fig 1a, and**
 297 **achieves better performance than reward learning baselines.** We show the normalized return
 298 in each AntMaze environment for all methods in Table 1. Without the need to learn an extra reward
 299 function, our method consistently achieves equivalent or better performance in each environment
 300 compared to other reward learning baselines. We observe that our method achieves comparable
 301 average performance to the context-agnostic method, given that goal sets are all very similar.⁹

302 **Reward model evaluation for reward learning baselines.** We also visualize the learned reward
 303 model from reward learning baselines¹⁰ to show how good they are at predicting the reward, and
 304 how it is related to the performance. Take “medium-diverse” and “large-diverse” environments as
 305 examples (see Figure 2, 3). For PDS, we can observe that the reward distribution for positive and
 306 negative samples are better separated in the large one than the medium one, explaining that it has
 307 better performance in the large-diverse environment than the medium-diverse one. Also, we observe
 308 that UDS+RP is consistently better at separating positive and negative distributions than plain RP, so
 309 we omit to compare with RP in the rest of the experiments. Intuitively, our method does not require
 310 reward learning thanks to the augmented MDP, which avoids the extra errors in reward prediction.

Env/Method	SDS (Ours)	PDS	RP	UDS+RP	Context-agnostic IQL
umaze	94.8±1.3	87.2±2.5	50.5±2.1	54.3±6.3	97.7±1.0
umaze diverse	72.8±7.7	73.2±3.1	72.8±2.6	71.5±4.3	65.5±10.5
medium play	75.8±1.9	35.2±8.2	0.5±0.3	0.3±0.3	75.2±3.4
medium diverse	84.5±5.2	3.8±1.7	0.5±0.5	0.8±0.5	76.0±3.7
large play	60.0±7.6	41.5±4.9	0±0	0±0	45.8±2.6
large diverse	36.8±6.9	28.8±6.3	0±0	0±0	46.7±5.4
average	70.8	45.0	20.7	21.2	67.8

Table 1: Normalized return in AntMaze-v2, averaged over 5 random seeds with standard errors.

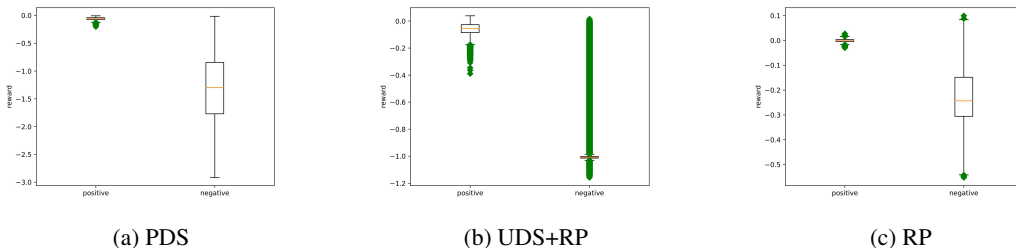


Figure 2: Reward model evaluation for the large-diverse environment. Green dots are outliers.

311 4.4 Modified AntMaze: Four Rooms

312 **Context-goal setup.** We partition the maze into four rooms and any state in the room would be a
 313 goal state. We use discrete room numbers (1,2,3,4) as contexts. As the agent always starts in Room
 314 1, the training and test context sets are Room 2,3,4. We use medium-play and large-play datasets.

⁹Also, we find that umaze is too easy such that even if the goal labeling is bad it still has a relatively high reward (since the maze is too small), so we also omit umaze in other experiments. Li et al. (2023b) show offline RL algorithms can learn good with goal-reaching data even when the rewards are wrong.

¹⁰We include the details for reward model evaluation in Appendix C.2.

315 **SDS achieves better performance than reward learning baselines under the setting in Figure 1b.** We show the normalized return (average success rate in percentage) in each modified Four Rooms environment for our method and baseline methods in Table 2, where our method consistently outperforms all baseline methods in each environment. We observe that the context agnostic method achieves rather high performance under this setting. This is because the number of rooms is only three, and the context agnostic method will learn to reach one room always with a high successful rate so the average is roughly 1/3, but it will not be the case in Section 4.5 when we have more test contexts. We also provide evaluation for reward learning in Figure 5.

Env/Method	SDS (Ours)	PDS	UDS+RP	Context-agnostic IQL
medium	78.2±1.2	26.3±1.6	14.0±0.9	32.6±0.8
large	73.3±1.9	14.0±2.7	21.6±21.3	28.1±0.3

Table 2: Average scores with standard errors over 5 random seeds from Four Rooms. The score for each run is the average success rate (%) of the other three rooms.

322
323

4.5 Modified AntMaze: Random Cells

324 **Context-goal setup.** We use the 2D locations as context but the distribution of the context is much more diverse than Section 4.3. For each maze map, we choose a set of non-wall 2D locations in the maze map, uniformly sample from it, and add uniform perturbations to get the training contexts. To construct the goal set given context, we obtain states with the 2D locations within the L_2 ball with a certain radius. For test distributions, we have two settings: 1) the same as the training distribution; 2) test contexts are drawn from a limited area that is far away from the starting point of the agent.

330 **SDS achieves better performance than reward learning baselines under the setting in Figure 1c.** We show the normalized return (average success rate in percentage) in each modified Random Cells environment for all methods in Table 3, where our method consistently outperforms all baseline methods in each environment, which also shows the generalization ability in the context space. We also provide reward visualization for reward learning baselines in Figure 6.

Env/Method	SDS (Ours)	PDS	UDS+RP	Context-agnostic IQL
medium	70.5±8.7	47.5±6.5	14.8±5.8	18.8±5.5
large	55.0±9.3	44.8±8.4	10.1±3.5	17.8±3.7

Table 3: Average scores with standard errors over 5 random seeds from Random Cells. The score for each run is the average success rate (%) of 5 random test contexts from the same training distribution.

335 **SDS also works with a different test context distribution.** We also test with a different distribution of random cells that are far away from the start with some specified threshold in each environment. We can observe that when tested with this different context distribution, SDS still consistently outperforms reward learning baselines.

Env/Method	SDS (Ours)	PDS	UDS+RP	Context-agnostic IQL
medium	63.8±11.9	31.5 ±18.0	2.2±0.9	4.3±1.7
large	62.6±6.4	44.6±7.6	1.1±0.6	0.8±0.8

Table 4: Average scores with standard errors over 5 random seeds from Random Cells. The score for each run is the average success rate (%) of 5 random test contexts of cells far away from the start.

339 5 Conclusion and Limitation

340 We propose a Simple Data Sharing technique for offline CGO problems. We prove SDS can learn near optimal policies so long as the offline data cover goal-reaching trajectories needed at the test time, without the need of negative labels. We also validate the efficacy of SDS experimentally, and we find it outperforms other reward-learning offline RL baselines across various CGO problem settings. We highlight SDS works under certain assumptions. As shown in our theoretical result in Section 3.3, the SDS technique would fail 1) if the dynamics dataset does not contain trajectories leading to the goal set of a given context, 2) the context-goal dataset does not cover the contexts and goals faced at test time, or 3) if the goal set does not cover reachable goals from initial states. While we believe SDS for its simplicity and theoretical guarantees would be useful in real-world settings (such as learning visual-language robot policies), our experimental setup is limited to low-dimensional simulation environments. Scaling up SDS empirically is an interesting future direction.

351 **References**

- 352 Marcin Andrychowicz, Filip Wolski, Alex Ray, Jonas Schneider, Rachel Fong, Peter Welinder, Bob
353 McGrew, Josh Tobin, Pieter Abbeel, and Wojciech Zaremba. Hindsight experience replay. In
354 *NeurIPS*, 2017.
- 355 André Barreto, Will Dabney, Rémi Munos, Jonathan J Hunt, Tom Schaul, Hado van Hasselt, and
356 David Silver. Successor features for transfer in reinforcement learning. In *NeurIPS*, 2017.
- 357 Yevgen Chebotar, Karol Hausman, Yao Lu, Ted Xiao, Dmitry Kalashnikov, Jacob Varley, Alex
358 Irpan, Benjamin Eysenbach, Ryan C Julian, Chelsea Finn, et al. Actionable models: Unsupervised
359 offline reinforcement learning of robotic skills. In *ICML*, 2021.
- 360 Ching-An Cheng, Tengyang Xie, Nan Jiang, and Alekh Agarwal. Adversarially trained actor critic
361 for offline reinforcement learning. In *ICML*, 2022.
- 362 Carlo D’Eramo, Davide Tateo, Andrea Bonarini, Marcello Restelli, and Jan Peters. Sharing knowl-
363 edge in multi-task deep reinforcement learning. In *ICLR*, 2020.
- 364 Justin Fu, Aviral Kumar, Ofir Nachum, George Tucker, and Sergey Levine. D4rl: Datasets for deep
365 data-driven reinforcement learning. *arXiv preprint arXiv:2004.07219*, 2020.
- 366 Scott Fujimoto and Shixiang Gu. A minimalist approach to offline reinforcement learning. In
367 *NeurIPS*, 2021.
- 368 Scott Fujimoto, David Meger, and Doina Precup. Off-policy deep reinforcement learning without
369 exploration. In *ICML*, 2019.
- 370 Meera Hahn, Devendra Singh Chaplot, Shubham Tulsiani, Mustafa Mukadam, James M Rehg, and
371 Abhinav Gupta. No rl, no simulation: Learning to navigate without navigating. In *NeurIPS*, 2021.
- 372 Assaf Hallak, Dotan Di Castro, and Shie Mannor. Contextual markov decision processes. *arXiv*
373 *preprint arXiv:1502.02259*, 2015.
- 374 Beining Han, Chongyi Zheng, Harris Chan, Keiran Paster, Michael R Zhang, and Jimmy Ba. Learn-
375 ing domain invariant representations in goal-conditioned block mdps. In *NeurIPS*, 2021.
- 376 Matteo Hessel, Hubert Soyer, Lasse Espeholt, Wojciech Czarnecki, Simon Schmitt, and Hado
377 Van Hasselt. Multi-task deep reinforcement learning with popart. In *AAAI*, 2019.
- 378 Hao Hu, Yiqin Yang, Qianchuan Zhao, and Chongjie Zhang. The provable benefit of unsupervised
379 data sharing for offline reinforcement learning. In *ICLR*, 2023.
- 380 Ying Jin, Zhuoran Yang, and Zhaoran Wang. Is pessimism provably efficient for offline rl? In *ICML*,
381 2021.
- 382 Leslie Pack Kaelbling. Learning to achieve goals. In *IJCAI*, 1993.
- 383 Dmitry Kalashnikov, Jacob Varley, Yevgen Chebotar, Benjamin Swanson, Rico Jonschkowski,
384 Chelsea Finn, Sergey Levine, and Karol Hausman. Mt-opt: Continuous multi-task robotic re-
385 inforcement learning at scale. *arXiv preprint arXiv:2104.08212*, 2021.
- 386 Ilya Kostrikov, Ashvin Nair, and Sergey Levine. Offline reinforcement learning with implicit q-
387 learning. In *ICLR*, 2021.
- 388 Aviral Kumar, Aurick Zhou, George Tucker, and Sergey Levine. Conservative q-learning for offline
389 reinforcement learning. In *NeurIPS*, 2020.
- 390 Alexander C Li, Lerrel Pinto, and Pieter Abbeel. Generalized hindsight for reinforcement learning.
391 In *NeurIPS*, 2020.
- 392 Anqi Li, Byron Boots, and Ching-An Cheng. Mahalo: Unifying offline reinforcement learning and
393 imitation learning from observations. In *ICML*, 2023a.
- 394 Anqi Li, Dipendra Misra, Andrey Kolobov, and Ching-An Cheng. Survival instinct in offline rein-
395 forcement learning. *arXiv preprint arXiv:2306.03286*, 2023b.

- 396 Corey Lynch, Mohi Khansari, Ted Xiao, Vikash Kumar, Jonathan Tompson, Sergey Levine, and
397 Pierre Sermanet. Learning latent plans from play. In *CORL*, 2020.
- 398 Yecheng Jason Ma, Jason Yan, Dinesh Jayaraman, and Osbert Bastani. Offline goal-conditioned
399 reinforcement learning via f -advantage regression. In *NeurIPS*, 2022.
- 400 Piotr Mirowski, Matthew Koichi Grimes, Mateusz Malinowski, Karl Moritz Hermann, Keith An-
401 derson, Denis Teplyashin, Karen Simonyan, Koray Kavukcuoglu, Andrew Zisserman, and Raia
402 Hadsell. Learning to navigate in cities without a map. In *NeurIPS*, 2018.
- 403 Dipendra K Misra, Jaeyong Sung, Kevin Lee, and Ashutosh Saxena. Tell me dave: Context-sensitive
404 grounding of natural language to manipulation instructions. *International Journal of Robotics*
405 *Research*, 35(1-3):281–300, 2016.
- 406 Ashvin Nair, Vitchyr Pong, Murtaza Dalal, Shikhar Bahl, Steven Lin, and Sergey Levine. Visual
407 reinforcement learning with imagined goals. In *NeurIPS*, 2018.
- 408 Suraj Nair and Chelsea Finn. Hierarchical foresight: Self-supervised learning of long-horizon tasks
409 via visual subgoal generation. In *ICLR*, 2019.
- 410 Tom Schaul, Daniel Horgan, Karol Gregor, and David Silver. Universal value function approxima-
411 tors. In *ICML*, 2015.
- 412 Avi Singh, Albert Yu, Jonathan Yang, Jesse Zhang, Aviral Kumar, and Sergey Levine. Cog:
413 Connecting new skills to past experience with offline reinforcement learning. *arXiv preprint*
414 *arXiv:2010.14500*, 2020.
- 415 Shagun Sodhani, Amy Zhang, and Joelle Pineau. Multi-task reinforcement learning with context-
416 based representations. In *ICML*, 2021.
- 417 Pei Sun, Henrik Kretzschmar, Xerxes Dotiwalla, Aurelien Chouard, Vijaysai Patnaik, Paul Tsui,
418 James Guo, Yin Zhou, Yuning Chai, Benjamin Caine, Vijay Vasudevan, Wei Han, Jiquan
419 Ngiam, Hang Zhao, Aleksei Timofeev, Scott Ettinger, Maxim Krivokon, Amy Gao, Aditya Joshi,
420 Yu Zhang, Jonathon Shlens, Zhifeng Chen, and Dragomir Anguelov. Scalability in perception for
421 autonomous driving: Waymo open dataset. In *CVPR*, 2020.
- 422 Yee Whye Teh, Victor Bapst, Wojciech Marian Czarnecki, John Quan, James Kirkpatrick, Raia
423 Hadsell, Nicolas Heess, and Razvan Pascanu. Distral: robust multitask reinforcement learning.
424 In *NeurIPS*, 2017.
- 425 Homer Rich Walke, Kevin Black, Tony Z. Zhao, Quan Vuong, Chongyi Zheng, Philippe Hansen-
426 Estruch, Andre Wang He, Vivek Myers, Moo Jin Kim, Max Du, Abraham Lee, Kuan Fang,
427 Chelsea Finn, and Sergey Levine. Bridgedata v2: A dataset for robot learning at scale. In *CORL*,
428 2023.
- 429 Benjamin Wilson, William Qi, Tanmay Agarwal, John Lambert, Jagjeet Singh, Siddhesh Khandel-
430 wal, Bowen Pan, Ratnesh Kumar, Andrew Hartnett, Jhony Kaesemodel Pontes, Deva Ramanan,
431 Peter Carr, and James Hays. Argoverse 2: Next generation datasets for self-driving perception
432 and forecasting. In *NeurIPS*, 2021.
- 433 Yifan Wu, George Tucker, and Ofir Nachum. Behavior regularized offline reinforcement learning.
434 *arXiv preprint arXiv:1911.11361*, 2019.
- 435 Tengyang Xie, Ching-An Cheng, Nan Jiang, Paul Mineiro, and Alekh Agarwal. Bellman-consistent
436 pessimism for offline reinforcement learning. In *NeurIPS*, 2021.
- 437 Rui Yang, Lin Yong, Xiaoteng Ma, Hao Hu, Chongjie Zhang, and Tong Zhang. What is essential
438 for unseen goal generalization of offline goal-conditioned rl? In *ICML*, 2023.
- 439 Albert Yu and Ray Mooney. Using both demonstrations and language instructions to efficiently learn
440 robotic tasks. In *ICLR*, 2023.
- 441 Tianhe Yu, Garrett Thomas, Lantao Yu, Stefano Ermon, James Zou, Sergey Levine, Chelsea Finn,
442 and Tengyu Ma. Mopo: model-based offline policy optimization. In *NeurIPS*, 2020.

- 443 Tianhe Yu, Aviral Kumar, Yevgen Chebotar, Karol Hausman, Sergey Levine, and Chelsea Finn.
444 Conservative data sharing for multi-task offline reinforcement learning. In *NeurIPS*, 2021.
- 445 Tianhe Yu, Aviral Kumar, Yevgen Chebotar, Karol Hausman, Chelsea Finn, and Sergey Levine.
446 How to leverage unlabeled data in offline reinforcement learning. In *ICML*, 2022.
- 447 Zhuangdi Zhu, Kaixiang Lin, Anil K Jain, and Jiayu Zhou. Transfer learning in deep reinforcement
448 learning: A survey. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2023.

449 A Detailed Related Work

450 **Goal-oriented RL** GO RL has been extensively studied (Kaelbling, 1993). Existing work focus
451 on two critical aspects of goal-oriented RL: (1) data relabeling and augmentation methods to make
452 better use of available data and (2) learning reusable skills to solve long-horizon problems by chain-
453 ing sub-goals or skills. For (1), hindsight relabeling methods (Andrychowicz et al., 2017; Li et al.,
454 2020) are effective in improving the learning efficiency of agents by reusing visited states in the
455 trajectories as successful goal examples. For (2), hierarchical methods for determining sub-goals,
456 and training goal reaching policies have been effective in long-horizon problems (Nair & Finn,
457 2019; Singh et al., 2020; Chebotar et al., 2021). Beyond data efficiency, another key objective of
458 goal-oriented RL is generalization, wherein a common representation of target goals is learned. Pop-
459 ular strategies for goal generalization include universal value function approximators (Schaul et al.,
460 2015), unsupervised representation learning (Nair et al., 2018; Nair & Finn, 2019; Han et al., 2021),
461 and pessimism induced generalization in offline GO formulations (Yang et al., 2023). Our CGO
462 framing enables both data reuse and goal generalization, by using rich contextual representations of
463 goals and a reduction to offline RL to combine dynamics and context-goal datasets.

464 **Offline RL** Offline RL methods have proven to be effective in GO problems as it also allows
465 learning a common set of sub-goals/skills (Chebotar et al., 2021; Ma et al., 2022; Yang et al., 2023).
466 A variety of approaches are used to mitigate the distribution shift between the collected datasets and
467 the trajectories likely to be generated by learnt policies: (1) constrain target policies to be close to the
468 dataset distribution (Fujimoto et al., 2019; Wu et al., 2019; Fujimoto & Gu, 2021), (2) incorporate
469 value pessimism for low-coverage or Out-Of-Distribution states and actions (Kumar et al., 2020; Yu
470 et al., 2020; Jin et al., 2021) and (3) adversarial training via a two-player game (Xie et al., 2021;
471 Cheng et al., 2022). Our SDS allows the use of generic offline RL algorithms to solve CGO problem
472 offline. We demonstrate its applicability with PSPI (Xie et al., 2021) and IQL (Kostrikov et al., 2021)
473 as our base offline RL algorithm in analyses (Section 3.3) and experiments (Section 4), respectively.

474 **Offline RL with unlabeled data** Our CGO setting is a special case of offline RL with unlabeled
475 data, or more broadly the offline policy learning from observations paradigm (Li et al., 2023a). There
476 only a subset of the offline data is labeled with rewards (in our setting, that is the contexts dataset, as
477 we don’t know which samples in the dynamics dataset are goals.). However, the MAHALO scheme
478 in (Li et al., 2023a) is much more general than necessary for CGO problems, and we show instead
479 that our simple data sharing scheme has better theoretical guarantees than MAHALO in Section 3.3.
480 In our experiments, we compare CGO with several offline RL algorithms designed for unlabeled
481 data: UDS (Yu et al., 2022) where unlabeled data is assigned zero rewards and PDS (Hu et al.,
482 2023) where a pessimistic reward function is estimated from a labeled dataset.

483 **Data-sharing in RL** Sharing information across multiple tasks is a promising approach to accel-
484 erate learning and to identify transferable features across tasks. In RL, both multi-task and transfer
485 learning settings have been studied under varying assumption on the shared properties and structures
486 of different tasks (Zhu et al., 2023; Teh et al., 2017; Barreto et al., 2017; D’Eramo et al., 2020). For
487 data sharing in CGO, we adopt the contextual MDP formulation (Hallak et al., 2015; Sodhani et al.,
488 2021), which enables knowledge transfer via high-level contextual cues. Prior work on offline RL
489 has also shown the utility of sharing data across tasks: hindsight relabeling and manual skill group-
490 ing (Kalashnikov et al., 2021), inverse RL (Li et al., 2020), sharing Q-value estimates (Yu et al.,
491 2021; Singh et al., 2020) and reward labeling (Yu et al., 2022; Hu et al., 2023).

492 B SDS +PSPI: Theoretical Analysis

493 In this section, we provide a detailed analysis for the instantiation of SDS using PSPI. We follow
494 the same notation for the value functions, augmented MDP and extended function classes as stated
495 in Section 2 and Section 3 in the main text.

496 B.1 Equivalence relations between original and Augmented MDP

497 We begin by showing that the optimal policy and any value function in the augmented MDP can
498 be expressed using their analogue in the original MDP. With the augmented MDP defined as $\overline{\mathcal{M}} :=$

499 $(\bar{\mathcal{X}}, \bar{\mathcal{A}}, \bar{R}, \bar{P}, \gamma)$ in Section 3.1, we first define the value function in the augmented MDP. For a policy
500 $\bar{\pi} : \bar{\mathcal{X}} \rightarrow \bar{\mathcal{A}}$, we define the Q function for the augmented MDP as

$$\bar{Q}^{\bar{\pi}}(x, a) := \mathbb{E}_{\bar{\pi}, \bar{P}} \left[\sum_{t=0}^{\infty} \gamma^t \bar{R}(x, a) \mid x_0 = x, a_0 = a \right]$$

501 Notice that we don't have a reaching time random variable T in this definition; instead the agent
502 would enter an absorbing state s^+ after taking a^+ in the augmented MDP. We can define similarly
503 $\bar{V}^{\bar{\pi}}(s) := \bar{Q}^{\bar{\pi}}(x, \bar{\pi})$.

504 **Remark B.1.** Let $\bar{Q}_R^{\bar{\pi}}$ be the extension of $Q^{\bar{\pi}}$ based on R . We have, for $x \notin G$, $\bar{Q}_R^{\bar{\pi}}(x, a) = \bar{Q}^{\bar{\pi}}(x, a)$
505 $\forall a \in \bar{\mathcal{A}}$, and for $x \in G$, $\bar{Q}_R^{\bar{\pi}}(x, a) = \bar{Q}^{\bar{\pi}}(x, a^+) = 1, \forall a \in \bar{\mathcal{A}}$.

506 By the construction of the augmented MDP, it is obvious that the following is true.

507 **Lemma B.2.** Given $\pi : \mathcal{X} \rightarrow \Delta(\mathcal{A})$, let $\bar{\pi}$ be its extension. For any $h : \mathcal{X} \times \mathcal{A} \rightarrow \mathbb{R}$, it holds

$$\mathbb{E}_{\pi, P} \left[\sum_{t=0}^T \gamma^t h(x, a) \right] = \mathbb{E}_{\bar{\pi}, \bar{P}} \left[\sum_{t=0}^{\infty} \gamma^t \tilde{h}^{\pi}(x, a) \mid x \notin \mathcal{X}^+ \right]$$

508 where T is the goal-reaching time (random variable) and we define $\tilde{h}^{\pi}(x, a^+) = h(x, \pi)$.

509 We can now relate the value functions between the two MDPs.

510 **Proposition B.3.** For a policy $\pi : \mathcal{S} \rightarrow \Delta(\mathcal{A})$, let $\bar{\pi}$ be its extension (defined above). We have for
511 all $x \in \mathcal{X}$, $a \in \mathcal{A}$,

$$\begin{aligned} Q^{\pi}(x, a) &\geq \bar{Q}^{\bar{\pi}}(x, a) \\ V^{\pi}(x) &= \bar{V}^{\bar{\pi}}(x) \end{aligned}$$

512 Conversely, for a policy $\xi : \bar{\mathcal{X}} \rightarrow \Delta(\bar{\mathcal{A}})$, define its restriction ξ on \mathcal{X} and \mathcal{A} by translating proba-
513 bility of ξ originally on a^+ to be uniform over \mathcal{A} . Then we have for all $s \in \mathcal{S}$, $a \in \mathcal{A}$

$$\begin{aligned} Q^{\xi}(x, a) &\geq \bar{Q}^{\xi}(x, a) \\ V^{\xi}(x) &\geq \bar{V}^{\xi}(x) \end{aligned}$$

514 *Proof.* The first direction follows from Lemma B.2. For the latter, whenever ξ takes a^+ at some
515 $x \notin G$, it has $\bar{V}^{\xi}(x) = 0$ but $V^{\xi}(x) \geq 0$ since there is no negative reward in the original MDP. By
516 performing a telescoping argument, we can derive the second claim. \square

517 By this lemma, we know the extension of π^* (i.e., $\bar{\pi}^*$) is also optimal to the augmented MDP and
518 $V^*(x) = \bar{V}^*(x)$ for $x \in \mathcal{X}$. Furthermore, we have a reduction that we can solve for the optimal
519 policy in the original MDP by the solving augmented MDP, since

$$V^{\xi}(d_0) - V^*(d_0) \leq V^{\xi}(d_0) - \bar{V}^*(d_0)$$

520 for all $\xi : \bar{\mathcal{X}} \rightarrow \Delta(\bar{\mathcal{A}})$. In particular,

$$\text{Regret}(\pi) := V^{\pi}(d_0) - V^*(d_0) = V^{\bar{\pi}}(d_0) - \bar{V}^*(d_0) =: \overline{\text{Regret}}(\bar{\pi}) \quad (5)$$

521 Since the augmented MDP replaces the random reaching time construction with an absorbing-state
522 version, the Q function $\bar{Q}^{\bar{\pi}}$ of the extended policy $\bar{\pi}$ satisfies the Bellman equation

$$\begin{aligned} \bar{Q}^{\bar{\pi}}(x, a) &= \bar{R}(x, a) + \gamma \mathbb{E}_{x' \sim \bar{P}(\cdot \mid x, a)} [\bar{Q}^{\bar{\pi}}(x', \bar{\pi})] \\ &=: \bar{\mathcal{T}}^{\bar{\pi}} \bar{Q}^{\bar{\pi}}(x, a) \end{aligned} \quad (6)$$

523 For $x \in \mathcal{X}$ and $a \in \mathcal{A}$, we show how the above equation can be rewritten in Q^{π} and R .

524 **Proposition B.4.** For $x \in \mathcal{X}$ and $a \in \mathcal{A}$,

$$\bar{Q}^{\bar{\pi}}(x, a) = 0 + \gamma \mathbb{E}_{x' \sim \bar{P}(\cdot \mid x, a)} [\max(R(x'), Q^{\pi}(x', \pi))]$$

525 For $a = a^+$, $\bar{Q}^{\bar{\pi}}(x, a^+) = \bar{R}(x, a^+) = R(x)$. For $x \in \mathcal{X}^+$, $\bar{Q}^{\bar{\pi}}(x, a) = 0$.

526 *Proof.* The proof follows from Lemma B.5 and the definition of \bar{P} . □

527 **Lemma B.5.** For $x \in \mathcal{X}$, $\bar{Q}^\pi(x, \bar{\pi}) = \max(R(x), Q^\pi(x, \pi))$

528 *Proof.* For $x \in \mathcal{X}$,

$$\begin{aligned}
\bar{Q}^\pi(x, \bar{\pi}) &= \begin{cases} \bar{Q}^\pi(x, a^+), & \text{if } x \in G \\ \bar{Q}^\pi(x, \pi), & \text{otherwise} \end{cases} && \text{(Because of definition of } \bar{\pi} \text{)} \\
&= \begin{cases} \bar{Q}^\pi(x, a^+), & \text{if } x \in G \\ Q^\pi(x, \pi), & \text{otherwise} \end{cases} && \text{(Because of Proposition B.3)} \\
&= \begin{cases} \bar{R}(x, a^+), & \text{if } x \in G \\ Q^\pi(x, \pi), & \text{otherwise} \end{cases} && \text{(Definition of augmented MDP)} \\
&= \begin{cases} R(x), & \text{if } x \in G \\ Q^\pi(x, \pi), & \text{otherwise} \end{cases} \\
&= \max(R(x), Q^\pi(x, \pi))
\end{aligned}$$

529 where in the last step we use $\bar{R}(x) = 1$ for $x \in G$ and $\bar{R}(x) = 0$ otherwise. □

530 B.2 Function Approximator Assumptions

531 In Theorem 3.1, we assume access to a policy class $\Pi = \{\pi : \mathcal{X} \rightarrow \Delta(\mathcal{A})\}$. We also assume access
532 to a function class $\mathcal{F} = \{f : \mathcal{X} \times \mathcal{A} \rightarrow [0, 1]\}$ and a function class $\mathcal{G} = \{g : \mathcal{X} \rightarrow [0, 1]\}$. We can
533 think of them as approximator for the Q function and the reward function of the original MDP.

534 Recall the zero-reward Bellman backup operator \mathcal{T}^π with respect to $P(s'|s, a)$ as defined in As-
535 sumption 3.3:

$$\mathcal{T}^\pi f(x, a) := \gamma \mathbb{E}_{x' \sim P_0(\cdot|x, a)}[f(x', \pi)]$$

536 where $P_0(x'|s, a) := P(s'|s, a)\mathbb{1}(c' = c)$. Note this definition is different from the one with
537 absorbing state s^+ in Section 2.3. Using this modified backup operator, we can show that the
538 following realizability assumption is true for the augmented MDP:

539 **Proposition B.6** (Realizability). *By Assumption 3.2 and Assumption 3.3, there is $f \in \mathcal{F}$ and $g \in \mathcal{G}$*
540 *such that $\bar{Q}^\pi = f_g$.*

541 *Proof.* By Assumption 3.3, there is $h \in \mathcal{F}$ such that $h(x, a) = \max(R(x), Q^\pi(x, a))$. By Proposi-
542 tion B.4, we have for $x \in \mathcal{X}$, $a \neq a^+$

$$\begin{aligned}
\bar{Q}^\pi(x, a) &= 0 + \gamma \mathbb{E}_{x' \sim \bar{P}(\cdot|x, a)}[\max(R(x'), Q^\pi(x', \pi))] \\
&= 0 + \gamma \mathbb{E}_{x' \sim P_0(\cdot|x, a)}[h(x, \pi)] \\
&= \mathcal{T}^\pi h \in \mathcal{F}
\end{aligned}$$

543 For $a = a^*$, we have $\bar{Q}^\pi(x, a^*) = \bar{R}(x, a^+) = R(x) \in \mathcal{G}$. Finally $\bar{Q}^\pi(x^+, a) = 0$ for $x^+ \in \mathcal{X}^+$.
544 Therefore, $\bar{Q}^\pi = f_g$ for some $f \in \mathcal{F}$ and $g \in \mathcal{G}$. □

545 B.3 Algorithm

546 In this section, we describe the instantiation of PSPI with SDS in detail along with the necessary
547 notation. As discussed in Section 3.3, our algorithm is based on the idea of reduction, which turns
548 the offline CGO problem into an standard offline RL problem in the augmented MDP. To this end,
549 we construct augmented datasets \bar{D}_{dyn} and \bar{D}_{goal} in Algorithm 1 as follows:

$$\begin{aligned}
\bar{D}_{\text{dyn}} &= \{(x_n, a_n, r_n, x'_n) | r_n = 0, x_n = (s_i, c_j), x'_n = (s'_i, c_j), a_n = a_i, (s_i, a_i, s'_i) \in D_{\text{dyn}}, (\cdot, c_j) \in D_{\text{goal}}\} \\
\bar{D}_{\text{goal}} &= \{(x_n, a^+, r_n, x_n^+) | r_n = 1, x_n = (s_n, c_n), x_n^+ = (s^+, c_n), (s_n, c_n) \in D_{\text{goal}}\}
\end{aligned}$$

550 For the analysis, we consider a simplified version of Algorithm 1 where we do not reuse the samples
551 in D_{dyn} . Specifically, for each sample $(s_i, a_i, s'_i) \in D_{\text{dyn}}$, we pair it with one sample $(\cdot, c_j) \in D_{\text{goal}}$
552 and do not reuse the sample from D_{dyn} . This can be naively done by pairing observed transitions and

553 context-goal pairs in both datasets when $|D_{\text{goal}}| \geq |D_{\text{dyn}}|$. In the analysis, we will state our results
 554 under this simplification.

555 With this construction, we have: $\bar{D}_{\text{dyn}} \sim \mu_{\text{dyn}}(s, a, s')\mu_{\text{goal}}(c)$ and $\bar{D}_{\text{goal}} \sim \mu_{\text{goal}}(c, s)\mathbb{1}(a =$
 556 $a^+)\mathbb{1}(s' = s^+)$. With abuse of notation, we write $\mu_{\text{dyn}}(x, a, x') = \mu_{\text{dyn}}(s, a, s')\mu_{\text{goal}}(c)$ and
 557 $\mu_{\text{goal}}(x, a, x') = \mu_{\text{goal}}(c, s)\mathbb{1}(a = a^+)\mathbb{1}(s' = s^+)$. Note that, $|\bar{D}_{\text{goal}}| = |D_{\text{goal}}|$ and $|\bar{D}_{\text{dyn}}| = |D_{\text{dyn}}|$
 558 as we are simply augmenting the observed states and actions without reusing samples. These two
 559 datasets have the standard tuple format, so we can run offline RL on $\bar{D}_{\text{dyn}} \cup \bar{D}_{\text{goal}}$.

560 **SDS +PSPI** We consider the information theoretic version of PSPI (Xie et al., 2021) which can
 561 be summarized as follows: For an MDP $(\mathcal{X}, \mathcal{A}, R, P, \gamma)$, given a tuple dataset $D = \{(x, a, r, x')\}$,
 562 a policy class Π , and a value class \mathcal{F} , it finds the policy through solving the two-player game:

$$\max_{\pi \in \Pi} \min_{f \in \mathcal{F}} f(d_0, \pi) \quad \text{s.t.} \quad \ell(f, f; \pi, D) - \min_{f' \in \mathcal{F}} \ell(f', f; \pi, D) \leq \epsilon_b \quad (7)$$

563 where $f(d_0, \pi) = \mathbb{E}_{x_0 \sim d_0}[f(x_0, \pi)]$, $\ell(f, f'; \pi, D) := \frac{1}{|D|} \sum_{(x, a, r, x') \in D} (f(x, a) - r - f'(x', \pi))^2$.
 564 The term $\ell(f, f; \pi, D) - \min_{f'} \ell(f', f; \pi, D)$ in the constraint is an empirical estimation of the
 565 Bellman error on f with respect to π on the data distribution μ , i.e. $\mathbb{E}_{x, a \sim \mu} [(f(x, a) - \mathcal{T}^\pi f(x, a))^2]$.
 566 It constrains the Bellman error to be small, since $\mathbb{E}_{x, a \sim \mu} [(Q^\pi(x, a) - \mathcal{T}^\pi Q^\pi(x, a))^2] = 0$.

567 Below we show how to run PSPI to solve the augmented MDP with offline dataset $\bar{D}_{\text{dyn}} \cup \bar{D}_{\text{goal}}$.
 568 To this end, we extend the policy class from Π to $\bar{\Pi}$, and the value class from \mathcal{F} to $\bar{\mathcal{F}}_{\mathcal{G}}$ using the
 569 function class \mathcal{G} based on the extensions defined in Section 3.1. One natural attempt is to implement
 570 equation 7 with the extended policy and value classes $\bar{\Pi}$ and $\bar{\mathcal{F}}$ and $\bar{D} = \bar{D}_{\text{dyn}} \cup \bar{D}_{\text{goal}}$. This would
 571 lead to the two player game:

$$\max_{\bar{\pi} \in \bar{\Pi}} \min_{\bar{f}_g \in \bar{\mathcal{F}}_{\mathcal{G}}} \bar{f}_g(d_0, \bar{\pi}) \quad \text{s.t.} \quad \ell(\bar{f}_g, \bar{f}_g; \bar{\pi}, \bar{D}) - \min_{\bar{f}'_g \in \bar{\mathcal{F}}_{\mathcal{G}}} \ell(\bar{f}'_g, \bar{f}_g; \bar{\pi}, \bar{D}) \leq \epsilon_b \quad (8)$$

572 However, equation 8 is not a well defined algorithm, because its usage of the extended policy $\bar{\pi}$ in
 573 the constraint requires knowledge of G , which is unknown to the agent.

574 Fortunately, we show that equation 8 can be slightly modified so that the implementation does not
 575 actually require knowing G . Here we use a property (Proposition B.4) that the Bellman equation of
 576 the augmented MDP:

$$\begin{aligned} \bar{Q}^{\bar{\pi}}(x, a) &= \bar{R}(x, a) + \gamma \mathbb{E}_{x' \sim \bar{P}(\cdot | x, a)} [\bar{Q}^{\bar{\pi}}(x', \bar{\pi})] \\ &= 0 + \gamma \mathbb{E}_{x' \sim \bar{P}(\cdot | x, a)} [\max(R(x'), Q^{\bar{\pi}}(x', \bar{\pi}))] \end{aligned}$$

577 for $x \in \mathcal{X}$ and $a \neq a^+$, and $\bar{Q}^{\bar{\pi}}(x, a) = 1$ for $x \in G$ and $a = a^+$.

578 We apply these two equalities to \bar{D}_{dyn} and \bar{D}_{goal} to construct our Bellman error estimates. Let
 579 $\phi(\bar{Q}^{\bar{\pi}}(x)) := \max(R(x), Q^{\bar{\pi}}(x, \bar{\pi}))$. We can rewrite the squared Bellman error on these two data
 580 distributions using the Bellman backup defined on the augmented MDP (see eq.6) as below:

$$\begin{aligned} \mathbb{E}_{x, a \sim \mu_{\text{dyn}}} [(\bar{Q}^{\bar{\pi}}(x, a) - \bar{\mathcal{T}}^{\bar{\pi}} \bar{Q}^{\bar{\pi}}(x, a))^2] &= \mathbb{E}_{x, a \sim \mu_{\text{dyn}}} [(\bar{Q}^{\bar{\pi}}(x, a) - 0 - \gamma \mathbb{E}_{x' \sim \bar{P}(\cdot | x, a)} [\phi(\bar{Q}^{\bar{\pi}})(x', \bar{\pi}))]^2] \\ \mathbb{E}_{x, a \sim \mu_{\text{goal}}} [(\bar{Q}^{\bar{\pi}}(x, a) - \bar{\mathcal{T}}^{\bar{\pi}} \bar{Q}^{\bar{\pi}}(x, a))^2] &= \mathbb{E}_{x, a \sim \mu_{\text{goal}}} [(\bar{Q}^{\bar{\pi}}(x, a^+) - 1)^2] \end{aligned}$$

582 We can construct an approximator $\bar{f}_g(x, a)$ for $\bar{Q}^{\bar{\pi}}(x, a)$. Substituting the estimator $\bar{f}_g(x, a)$ for
 583 $\bar{Q}^{\bar{\pi}}(x, a)$ in the squared Bellman errors above and approximating them by finite samples, we derive
 584 the empirical losses below.

$$\ell_{\text{dyn}}(\bar{f}_g, \bar{f}'_g; \bar{\pi}) := \frac{1}{|\bar{D}_{\text{dyn}}|} \sum_{(x, a, r, x') \in \bar{D}_{\text{dyn}}} (f(x, a) - \gamma \max(g'(x'), f'(x', \bar{\pi})))^2 \quad (9)$$

$$\ell_{\text{goal}}(\bar{f}_g) := \frac{1}{|\bar{D}_{\text{goal}}|} \sum_{(x, a, r, x') \in \bar{D}_{\text{goal}}} (g(x) - 1)^2 \quad (10)$$

585 where we use $\phi(\bar{f}_g)(x, a) = \max(g(x), f(x, a))$ and for $x \notin \mathcal{X}^+$, $\bar{f}_g(x, a) = f(x, a)\mathbb{1}(a \neq$
 586 $a^+) + g(x)\mathbb{1}(a = a^+)$.

587 Using this loss, we define the two-player game of PSPI for the augmented MDP:

$$\begin{aligned}
& \max_{\pi \in \Pi} \min_{\bar{f}_g \in \bar{\mathcal{F}}} \bar{f}_g(d_0, \bar{\pi}) \\
\text{s.t. } & \ell_{\text{dyn}}(\bar{f}_g, \bar{f}_g; \bar{\pi}) - \min_{\bar{f}'_{g'} \in \bar{\mathcal{F}}} \ell_{\text{dyn}}(\bar{f}'_{g'}, \bar{f}_g; \bar{\pi}) \leq \epsilon_{\text{dyn}} \\
& \ell_{\text{goal}}(\bar{f}_g) \leq 0
\end{aligned} \tag{11}$$

588 Notice $\bar{f}_g(d_0, \bar{\pi}) = f(d_0, \pi)$. Therefore, this problem can be solved using samples from D without
589 knowing G .

590 B.4 Analysis

591 **Covering number** We first define the covering number on the function classes \mathcal{F} , \mathcal{G} , and Π ¹¹. For
592 \mathcal{F} and \mathcal{G} , we use the L_∞ metric. We use $\mathcal{N}_\infty(\mathcal{F}, \epsilon)$ and $\mathcal{N}_\infty(\mathcal{G}, \epsilon)$ to denote their ϵ -covering
593 numbers. For Π , we use the L_∞ - L_1 metric, i.e., $\|\pi_1 - \pi_2\|_{\infty, 1} := \sup_{x \in \mathcal{X}} \|\pi_1(\cdot|x) - \pi_2(\cdot|x)\|_1$. We
594 use $\mathcal{N}_{\infty, 1}(\Pi, \epsilon)$ to denote its ϵ -covering number.

595 **High-probability Events** First, we show $\bar{Q}^{\bar{\pi}}$ has small empirical errors.

596 **Lemma B.7.** *With probability at least $1 - \delta$, it holds for all $\pi \in \Pi$,*

$$\begin{aligned}
& \ell_{\text{dyn}}(\bar{Q}^{\bar{\pi}}, \bar{Q}^{\bar{\pi}}; \bar{\pi}) - \min_{\bar{f}'_{g'} \in \bar{\mathcal{F}}} \ell_{\text{dyn}}(\bar{f}'_{g'}, \bar{Q}^{\bar{\pi}}; \bar{\pi}) \leq \epsilon_{\text{dyn}} \\
& \ell_{\text{goal}}(\bar{Q}^{\bar{\pi}}) \leq 0
\end{aligned}$$

597 where¹²

$$\epsilon_{\text{dyn}} = O\left(\frac{\log\left(\mathcal{N}_\infty\left(\mathcal{F}, \frac{1}{|D_{\text{dyn}}|}\right)\mathcal{N}_\infty\left(\mathcal{G}, \frac{1}{|D_{\text{dyn}}|}\right)\mathcal{N}_{\infty, 1}\left(\Pi, \frac{1}{|D_{\text{dyn}}|}\right)/\delta\right)}{|D_{\text{dyn}}|}\right)$$

598 *Proof.* Note $\bar{Q}^{\bar{\pi}} = \bar{f}_g$ for some $f \in \mathcal{F}$ and $g \in \mathcal{G}$ (Proposition B.6) and

$$0 = \mathbb{E}_{x, a \sim \mu_{\text{dyn}}}[(\bar{Q}^{\bar{\pi}}(x, a) - \bar{T}^{\bar{\pi}}\bar{Q}^{\bar{\pi}}(x, a))^2] = \mathbb{E}_{x, a \sim \mu_{\text{dyn}}}[(\bar{Q}^{\bar{\pi}}(x, a) - 0 - \gamma\mathbb{E}_{x' \sim \bar{P}(\cdot|x, a)}[\phi(\bar{Q}^{\bar{\pi}})(x', \pi)])^2]$$

599 Following a similar proof of Theorem 8 of (Cheng et al., 2022), we can derive ϵ_{dyn} . On the other
600 hand, $\ell_{\text{goal}}(\bar{f}_g) = 0$ because the reward $R(x)$ is deterministic. \square

601 Next, we show that with high probability the empirical error can upper bound the population error.

602 **Lemma B.8.** *For all $f \in \mathcal{F}, g \in \mathcal{G}$ satisfying*

$$\begin{aligned}
& \ell_{\text{dyn}}(\bar{f}_g, \bar{f}_g; \bar{\pi}) - \min_{\bar{f}'_{g'} \in \bar{\mathcal{F}}} \ell_{\text{dyn}}(\bar{f}'_{g'}, \bar{f}_g; \bar{\pi}) \leq \epsilon_{\text{dyn}} \\
& \ell_{\text{goal}}(\bar{f}_g) \leq 0
\end{aligned}$$

603 *With probability at least $1 - \delta$, for any $f \in \mathcal{F}, g \in \mathcal{G}$*

$$\begin{aligned}
& \|\bar{f}_g(x, a) - \gamma\mathbb{E}_{x' \sim \bar{P}(\cdot|x, a)}[\max(g(x'), f(x', \pi))]\|_{\mu_{\text{dyn}}} \leq O(\sqrt{\epsilon_{\text{dyn}}}) \\
& \|g(x) - 1\|_{\mu_{\text{goal}}} \leq O\left(\sqrt{\frac{\log\left(\frac{\mathcal{N}_\infty\left(\mathcal{G}, \frac{1}{|D_{\text{goal}}|}\right)}{\delta}\right)}{|D_{\text{goal}}|}}\right) =: \sqrt{\epsilon_{\text{goal}}}
\end{aligned}$$

604 *Proof.* This follows from Theorem 9 of (Cheng et al., 2022). \square

¹¹For finite function classes, the resulting performance guarantee will depend on $|\mathcal{F}|, |\mathcal{G}|$ and $|\Pi|$ instead of the covering numbers as stated in Theorem 3.1.

¹²Technically, we can remove $\mathcal{N}_\infty\left(\mathcal{G}, \frac{1}{|D_{\text{dyn}}|}\right)$ in the upper bound, but we include it here for a cleaner presentation.

605 **Pessimistic Estimate** We show the empirical value estimate found in equation 11 is pessimistic.

606 **Lemma B.9.** Given π , let \bar{f}_g^π denote the minimizer in equation 11. With high probability,
 607 $\bar{f}_g^\pi(d_0, \bar{\pi}) \leq Q^\pi(d_0, \pi)$

608 *Proof.* By Lemma B.7, we have $\bar{f}_g^\pi(d_0, \bar{\pi}) \leq \bar{Q}_R^\pi(d_0, \bar{\pi}) = Q^\pi(d_0, \pi)$. □

609 Next we bound the amount of underestimation.

610 **Lemma B.10.** Suppose $x_0 \sim d_0$ is not in G almost surely. For any $\pi \in \Pi$,

$$\begin{aligned} & Q^\pi(d_0, \pi) - \bar{f}_g^\pi(d_0, \bar{\pi}) \\ & \leq \mathbb{E}_\pi \left[\sum_{t=0}^{T-1} \gamma^t (\gamma \max(g^\pi(x_{t+1}), f^\pi(x_{t+1}, \pi)) - f^\pi(x_t, a_t)) + \gamma^T (R(x_T) - g^\pi(x_T)) \right] \end{aligned}$$

611 Note that in a trajectory $x_T \in G$ whereas $x_t \notin G$ for $t < T$ by definition of T .

612 *Proof.* Let $\bar{f}_g^\pi = (f^\pi, g^\pi)$ be the empirical minimizer. By performance difference lemma, we can
 613 write

$$\begin{aligned} & (1 - \gamma)Q^\pi(d_0, \pi) - (1 - \gamma)\bar{f}_g^\pi(d_0, \bar{\pi}) \\ & = (1 - \gamma)\bar{Q}^\pi(d_0, \bar{\pi}) - (1 - \gamma)\bar{f}_g^\pi(d_0, \bar{\pi}) \\ & = \mathbb{E}_{\bar{d}^\pi} [\bar{R}(x, a) + \gamma \bar{f}_g^\pi(x', \bar{\pi}) - \bar{f}_g^\pi(x, a)] \end{aligned}$$

614 where with abuse of notation we define $\bar{d}^\pi(x, a, x') := \bar{d}^\pi(x, a)\bar{P}(x'|x, a)$, where $\bar{d}^\pi(x, a)$ is the
 615 average state-action distribution of $\bar{\pi}$ in the augmented MDP.

616 In the above expectation, for $x \in G$, we have $a = a^+$ and $x^+ = (s^+, c)$ after taking a^+ at $x = (s, c)$,
 617 which leads to

$$\bar{R}(x, a) + \gamma \bar{f}_g^\pi(x', \bar{\pi}) - \bar{f}_g^\pi(x, a) = \bar{R}(x, a^+) + \gamma \bar{f}_g^\pi(x^+, \bar{\pi}) - \bar{f}_g^\pi(x, a^+) = R(x) - g^\pi(x)$$

618 For $x \notin G$ and $x \notin \mathcal{X}^+$, we have $a \neq a^+$ and $x' \notin \mathcal{X}^+$; therefore

$$\begin{aligned} \bar{R}(x, a) + \gamma \bar{f}_g^\pi(x', \bar{\pi}) - \bar{f}_g^\pi(x, a) & = R(x) + \gamma \bar{f}_g^\pi(x', \bar{\pi}) - f^\pi(x, a) \\ & \leq \gamma \max(g^\pi(x'), f^\pi(x', \pi)) - f^\pi(x, a) \end{aligned}$$

619 where the last step is because of the definition of \bar{f}_g^π . For $x \in \mathcal{X}^+$, we have $x \in \mathcal{X}^+$ and the reward
 620 is zero, so

$$\bar{R}(x, a) + \gamma \bar{f}_g^\pi(x', \bar{\pi}) - \bar{f}_g^\pi(x, a) = 0$$

621 Therefore, we can derive

$$\begin{aligned} & (1 - \gamma)Q^\pi(x_0, \pi) - (1 - \gamma)\bar{f}_g^\pi(x_0, \bar{\pi}) \\ & \leq \mathbb{E}_{\bar{d}^\pi} [\gamma \max(g^\pi(x'), f^\pi(x', \pi)) - f^\pi(x, a) | x \notin G, x \notin \mathcal{X}^+] + \mathbb{E}_{\bar{d}^\pi} [R(x) - g^\pi(x) | x \in G] \end{aligned}$$

622 Finally, using Lemma B.2 we can have the final upper bound.

623 □

624 **B.5 Main Result: Performance Bound**

625 Let π^\dagger be the learned policy and let $\bar{f}_g^{\pi^\dagger}$ be the learned function approximators. For any comparator
 626 policy π , let $\bar{f}_g^\pi = (f^\pi, g^\pi)$ be the estimator of π on the data. We have.

$$\begin{aligned}
 & V^\pi(d_0) - V^{\pi^\dagger}(d_0) \\
 &= Q^\pi(d_0, \pi) - Q^{\pi^\dagger}(d_0, \pi^\dagger) \\
 &= Q^\pi(d_0, \pi) - \bar{f}_g^{\pi^\dagger}(d_0, \bar{\pi}^\dagger) + \bar{f}_g^{\pi^\dagger}(d_0, \bar{\pi}^\dagger) - Q^{\pi^\dagger}(d_0, \pi^\dagger) \\
 &\leq Q^\pi(d_0, \pi) - \bar{f}_g^{\pi^\dagger}(d_0, \bar{\pi}^\dagger) \\
 &\leq Q^\pi(d_0, \pi) - \bar{f}_g^\pi(d_0, \bar{\pi}) \\
 &\leq \mathbb{E}_{\pi, P} \left[\sum_{t=0}^{T-1} \gamma^t (\gamma \max(g^\pi(x_{t+1}), f^\pi(x_{t+1}, \pi)) - f^\pi(x_t, a_t)) + \gamma^T (R(x_T) - g^\pi(x_T)) \right] \\
 &\leq \mathbb{E}_{\pi, P} \left[\sum_{t=0}^{T-1} \gamma^t |\gamma \max(g^\pi(x_{t+1}), f^\pi(x_{t+1}, \pi)) - f^\pi(x_t, a_t)| + \gamma^T |R(x_T) - g^\pi(x_T)| \right] \\
 &\leq \mathfrak{C}_{\text{dyn}}(\pi) \mathbb{E}_{\mu_{\text{dyn}}} [|\gamma \max(g^\pi(x'), f^\pi(x', \pi)) - f^\pi(x, a)|] + \mathfrak{C}_{\text{goal}}(\pi) \mathbb{E}_{\mu_{\text{goal}}} [|g(x) - 1|] \\
 &\leq \mathfrak{C}_{\text{dyn}}(\pi) \sqrt{\epsilon_{\text{dyn}}} + \mathfrak{C}_{\text{goal}}(\pi) \sqrt{\epsilon_{\text{goal}}}
 \end{aligned}$$

627 where $\mathfrak{C}_{\text{dyn}}(\pi)$ and $\mathfrak{C}_{\text{goal}}(\pi)$ are the concentrability coefficients defined in Definition 3.4.

628 **Theorem B.11.** Let π^\dagger denote the learned policy of SDS + PSPI with datasets D_{dyn} and D_{goal} ,
 629 using value function classes $\mathcal{F} = \{\mathcal{X} \times \mathcal{A} \rightarrow [0, 1]\}$ and $\mathcal{G} = \{\mathcal{X} \rightarrow [0, 1]\}$. Under realizability
 630 and completeness assumptions as stated in Assumption 3.2 and Assumption 3.3 respectively, with
 631 probability $1 - \delta$, it holds, for any $\pi \in \Pi$,

$$J(\pi) - J(\pi^\dagger) \leq \mathfrak{C}_{\text{dyn}}(\pi) \sqrt{\epsilon_{\text{dyn}}} + \mathfrak{C}_{\text{goal}}(\pi) \sqrt{\epsilon_{\text{goal}}}$$

632 where

$$\epsilon_{\text{dyn}} = O \left(\frac{\log \left(\mathcal{N}_\infty \left(\mathcal{F}, \frac{1}{|D_{\text{dyn}}|} \right) \mathcal{N}_\infty \left(\mathcal{G}, \frac{1}{|D_{\text{dyn}}|} \right) \mathcal{N}_{\infty, 1} \left(\Pi, \frac{1}{|D_{\text{dyn}}|} \right) / \delta \right)}{|D_{\text{dyn}}|} \right),$$

633 and,

$$\epsilon_{\text{goal}} = O \left(\frac{\log \left(\mathcal{N}_\infty \left(\mathcal{G}, \frac{1}{|D_{\text{goal}}|} \right) / \delta \right)}{|D_{\text{goal}}|} \right)$$

634 are statistical errors, and $\mathfrak{C}_{\text{dyn}}(\pi)$ and $\mathfrak{C}_{\text{goal}}(\pi)$ are concentrability coefficients which decrease as
 635 the data coverage increases.

636 **C Experimental details**

637 **C.1 Hyperparameters and experimental settings**

638 **IQL.** For IQL, we keep the hyperparameter of $\gamma = 0.99$, $\tau = 0.9$, $\beta = 10.0$, and $\alpha = 0.005$
 639 in Kostrikov et al. (2021), and tune other hyperparameters on the antmaze-medium-play-v2 envi-
 640 ronment and choose batch size = 1024 from candidate choices $\{256, 512, 1024, 2046\}$, learning rate
 641 = 10^{-4} from candidate choices $\{5 \cdot 10^{-5}, 10^{-4}, 3 \cdot 10^{-4}\}$ and 3 layer MLP with RuLU activating
 642 and 256 hidden units for all networks. We use the same set of IQL hyperparameters for both our
 643 methods and all the baseline methods included in Section 4.2, and apply it to all environments.

644 **RP.** For naive reward prediction, we use the full context-goal dataset as positive data, and train
 645 a reward model with 3-layer MLP and ReLU activations, learning rate = 10^{-4} , batch size = 1024,
 646 and training for 100 epochs for convergence. To label the transition dataset, we need to find some
 647 appropriate threshold to label states predicted as goals given contexts. We choose the percentile as

648 5% in the reward distribution evaluated by the context-goal set as the threshold to label goals in the
 649 antmaze-medium-play-v2 environment, from candidate choices {0%, 5%, 10%}. Then we apply it
 650 to all environments. Another trick we apply for the reward prediction is that instead of predicting 0
 651 for the context-goal dataset, we let it predict 1 but shift the reward prediction by -1 during reward
 652 evaluation, which prevents the model from learning all 0 weights. Similar tricks are also used in
 653 other reward learning baselines.

654 **UDS+RP.** We use the same structure and training procedure for the reward model as RP, except
 655 that we also randomly sample a minibatch of “negative” contextual transitions with the same batch
 656 size for a balanced distribution, which is constructed by randomly sampling combinations of a state
 657 in the trajectory-only dataset and a context from the context-goal dataset. To create a balanced
 658 distribution of positive and negative samples, we sample from each dataset with equal probability.
 659 For the threshold, we choose the percentile as 5% in the reward distribution evaluated by the context-
 660 goal set as the threshold to label goals in the antmaze-medium-play-v2 environment, from candidate
 661 choices {0%, 5%, 10%}. Then we apply it to all environments.

662 **PDS.** We use the same structure and training procedure for the reward model as RP, except that
 663 we train an ensemble of 10 networks as in Hu et al. (2023). To select the threshold percentile and
 664 the pessimistic weight k , we choose the percentile as 0% in the reward distribution evaluated by the
 665 context-goal set as the threshold to label goals from candidate choices {0%, 5%, 10%}, and $k = 15$
 666 from the candidate choices {5,10,15,20} in the antmaze-medium-play-v2 environment. Then we
 667 apply them to all environments.

668 **SDS (ours).** We do not require extra parameters other than the possibility of sampling from the
 669 real and fake transitions. Intuitively, we should sample from both datasets with the same probability
 670 to create an overall balanced distribution. Empirically, we also find that the balance distribution
 671 generates the best result.

672 C.2 Reward model evaluation

673 For reward learning baselines, we evaluate the learned reward model: we construct the positive
 674 dataset from context-goal examples, and the negative dataset from the combination of the context
 675 set and all states in the trajectory-only data, using the oracle context-goal function defined in the
 676 environment to filter out positive ones. We then evaluate the predicted reward on both positive and
 677 negative datasets, generating boxplots to visualize the distributions of the predicted reward for both
 678 datasets. The purpose of the reward model evaluation is to showcase whether the learned reward
 679 function can successfully capture context-goal relationships.

680 D More reward model evaluations

681 Here we present boxplots for reward models with experimental setups in Section 4.3, 4.4 and 4.5.

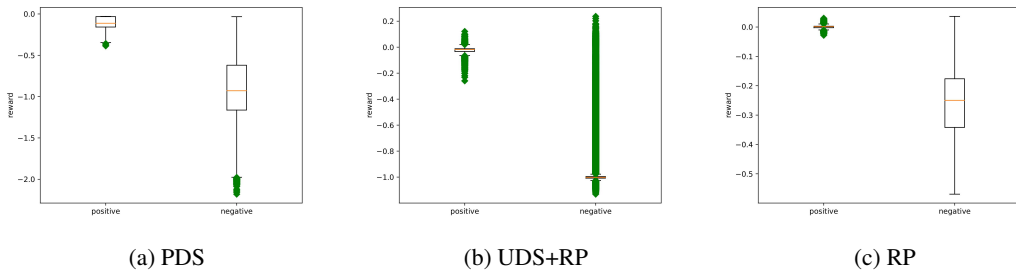


Figure 3: Reward model evaluation for the medium-diverse environment in Section 4.3. Green dots are outliers.

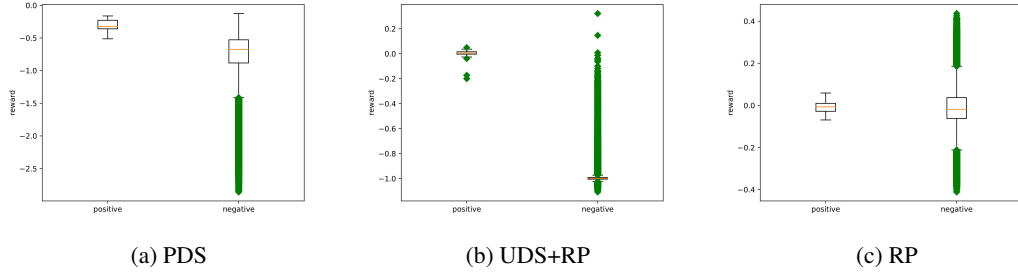


Figure 4: Reward model evaluation for the umaze-diverse environment in Section 4.3. Green dots are outliers.

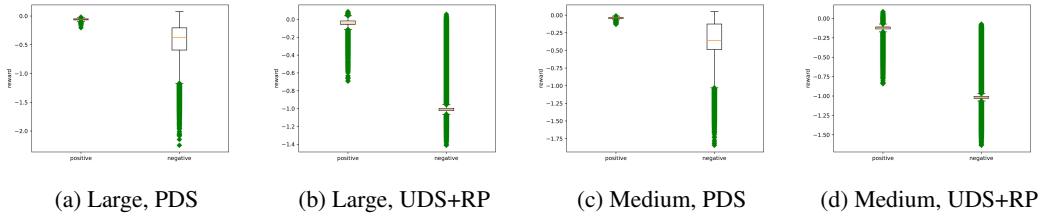


Figure 5: Reward model evaluation for Four Rooms in Section 4.4. Green dots are outliers.

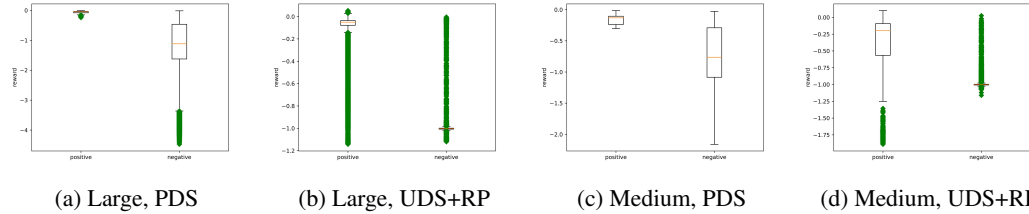


Figure 6: Reward evaluation for Random Cells in Section 4.5 (the test context distribution is the same as training). Green dots are outliers.

682 E Adding out-of-distribution (OOD) goal examples in the context-goal set

683 We include another table with a slightly different setting compared with Section 4.4: for each goal
 684 set given the context in the training context-goal set, we add some extra random states that are out
 685 of the original range of the state space as out-of-distribution goal examples (which are not covered
 686 by the trajectory-only dataset). The results are shown in Table 5, which is similar to the results in
 687 Section 4.4, showing that these methods are robust to extra OOD goal examples.

Env/Method	Ours	PDS	UDS+RP
medium	78.9±1.6	23.5±1.2	13.4±1.2
large	70.0±5.7	9.0±2.6	22.5±0.9

Table 5: Average scores with standard errors over 5 random seeds from Four Rooms, with extra OOD goal examples in the context-goal dataset. The reported score is the average success rate of three rooms, and the evaluation of each room requires 100 episodes.