
Safe Decision Transformer with Learning-based Constraints

Ruhan Wang
Indiana University
ruhwang@iu.edu

Dongruo Zhou
Indiana University
dz13@iu.edu

Abstract

1 In the field of safe offline reinforcement learning (RL), the objective is to utilize
2 offline data to train a policy that maximizes long-term rewards while adhering to
3 safety constraints. Recent work, such as the Constrained Decision Transformer
4 (CDT) [30], has utilized the Transformer [38] architecture to build a safe RL
5 agent that is capable of dynamically adjusting the balance between safety and task
6 rewards. However, it often lacks the stitching ability to output policies that are
7 better than those existing in the offline dataset, similar to other Transformer-based
8 RL agents like the Decision Transformer (DT) [7]. We introduce the Constrained
9 Q-learning Decision Transformer (CQDT) to address this issue. At the core of
10 our approach is a novel trajectory relabeling scheme that utilizes learned value
11 functions, with careful consideration of the trade-off between safety and cumula-
12 tive rewards. Experimental results show that our proposed algorithm outperforms
13 several baselines across a variety of safe offline RL benchmarks.

14 1 Introduction

15 Recent studies have demonstrated the Transformer’s [38] state-of-the-art performance across a range
16 of applications, including natural language processing [40, 4] and computer vision [27, 3]. When
17 applied to the domain of Reinforcement Learning (RL), a recent trend is to treat the decision-making
18 problem as a sequence modeling problem, using auto regressive models such as Transformer which
19 maps the history information directly to the action [7] or the next state [18]. Notably, the Decision
20 Transformer (DT) [7] effectively extends the Transformer architecture to offline RL tasks, showcasing
21 strong performance, particularly in sequential modeling. However, it is worth noting that while DT
22 excels in maximizing long-term rewards, it may not always align with the complexities of real-world
23 tasks. In practice, many tasks cannot be simplified solely into optimizing a single scalar reward
24 function, and the presence of various constraints significantly narrows the spectrum of feasible
25 solutions [12]. Such a setting is called safe RL, which has been studied in lots of safety-related
26 decision-making problems. For instance, it is crucial that robots interacting with humans in human-
27 machine environments prioritize human safety and avoid causing harm. In the realm of recommender
28 systems, it is imperative to avoid recommending false or racially discriminatory information to users.
29 Similarly, when self-driving cars operate in real-world environments, ensuring safety is paramount
30 [36, 14, 31].

31 Constrained Decision Transformer (CDT) [30] serves as a pioneering work which extends the
32 Transformer-based RL to the safe RL regime, which builds upon the foundation of the DT while
33 incorporating essential safety constraints. CDT’s core objective is to acquire a safe policy from an
34 offline dataset. Its distinguishing feature is the ability to maintain zero-shot adaptability across a
35 range of constraint thresholds, rendering it highly suitable for real-world reinforcement learning
36 applications burdened by constraints. While CDT demonstrates exceptional competitiveness in safe
37 offline reinforcement learning tasks, it shares a common limitation with DT—an absence of the
38 ‘stitching’ capability. This crucial ability involves integrating sub-optimal trajectory segments to form

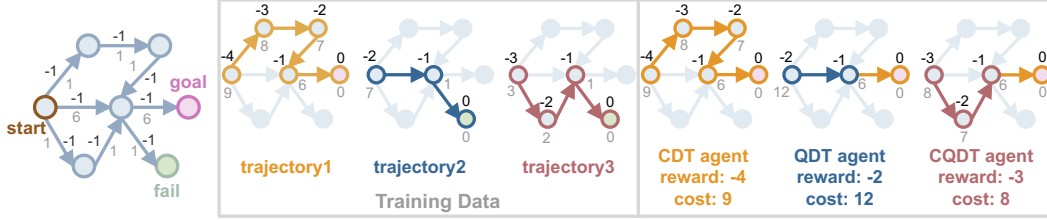


Figure 1: The toy example is presented where arrows denote the node connections. Here we use black number to denote rewards and gray nodes to denote costs. The demonstration example shows that cost return-based method (CDT) fails to find the shortest path to the goal since it lacks the stitching ability. In contrast, Q-learning-based DT (QDT) finds the shortest path, while it violates the safety constraints. Our proposed CQDT enjoys the superiority of both methods.

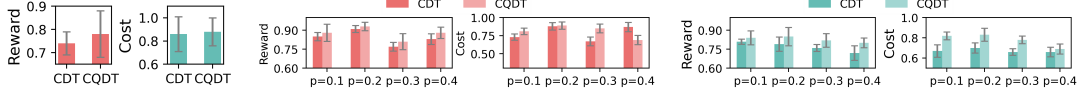


Figure 2: Performance comparison between CDT and CQDT on reward-suboptimal and cost-suboptimal datasets with results averaged over three random seeds. **Figure 3:** Stitching ability comparison between CDT and CQDT on cost-suboptimal datasets with results averaged over three random seeds. **Figure 4:** Stitching ability comparison between CDT and CQDT on reward-suboptimal datasets with results averaged over three random seeds.

39 a near-optimal trajectory, a pivotal characteristic highly desired in offline reinforcement learning
 40 agents. What is more challenging is that for RL with safety constraints, the agent not only needs to
 41 stitch trajectories to achieve better reward feedback but also needs to guarantee that the obtained
 42 policy is still *safe* after stitching, not being affected by unsafe trajectories in the offline dataset. Thus,
 43 we raise the following question:

44 **Can we design DT-based algorithms that output safe policies with the stitching ability?**

45 We answer this question affirmatively. To better demonstrate our algorithm design, we propose a toy
 46 example involving finding the path that maximizes reward under cost constraints, as illustrated in
 47 Figure 1. The task’s objective is to identify a path with the highest reward while adhering to a cost
 48 constraint (cost limitation = 10). The training data covers segments of the optimal path, but none
 49 of the training data trajectories encompass the entire optimal path. The agent must synthesize these
 50 fragmented trajectories to determine the optimal path within the cost constraint. For the existing
 51 DT-based RL algorithms such as Q-learning Decision Transformer (QDT) [42], they are able to
 52 stitch suboptimal trajectories into the optimal ones, but they are unable to maintain safety during
 53 the stitching phase. For the existing DT-based safe RL methods such as CDT, they can only obtain
 54 suboptimal policies while satisfying the safety guarantee, due to the lack of stitching ability. Thus,
 55 we propose the **Constrained Q-learning Decision Transformer (CQDT)** method to address the
 56 issues in existing methods, as shown in Figure 1. The main contributions are listed as follows.

- 57 • CQDT seeks to improve the quality of the training dataset by utilizing cost and reward critic
- 58 functions from the Constraints Penalized Q-learning framework [41] to relabel the return-to-go
- 59 (RTG) and cost-to-go (CTG) values in the training trajectories. This relabeled dataset is subsequently
- 60 used to train a Decision Transformer-based safe RL method, such as CDT.
- 61 • We provide a comparative analysis of our CQDT against various safe RL benchmarks across
- 62 different RL task settings. The results are summarized in Figure 2, demonstrate that CQDT
- 63 consistently outperforms all existing benchmarks in several offline safe RL tasks.
- 64 • We also show that our proposed CQDT enjoys better stitching ability compared with CDT, which
- 65 suggests that CQDT can better utilize suboptimal trajectories in the offline dataset. The results are
- 66 summarized in Figure 3 and Figure 4.

67 **2 Related Work**

68 **Offline Reinforcement Learning.** Offline Reinforcement Learning refers to a data-driven approach
 69 to the classic RL problem [25]. It focuses on deriving policies from pre-existing data without

70 requiring further interaction with the environment. The practical applications of offline RL are
 71 extensive, spanning domains such as healthcare [13] and the acquisition of robotic manipulation
 72 skills [19]. However, its inherent challenge lies in the potential disparity between the learned policy
 73 and the behavior that generated the initial offline data [21]. Addressing this challenge has spurred
 74 the development of various methodologies, including constraining the learned policy close from the
 75 behavior policy [11, 21, 22, 44, 17, 20, 39, 32, 18]. There is a recent line of work aiming at providing
 76 a Transformer-based policy without explicitly constraining the distribution shift issue [7]. Our work
 77 falls into that regime, which uses a Transformer as the policy model focusing on the safe RL setting.

78 **Offline Safe Reinforcement Learning.** Offline Safe Reinforcement Learning aims to acquire con-
 79 strained policies from pre-collected datasets, ensuring adherence to safety requirements throughout
 80 the learning process. This approach amalgamates techniques from both offline RL and safe RL, lever-
 81 aging methodologies from both domains [23]. Certain methods tackle the constrained optimization
 82 problem through stationary distribution correction, employing Lagrangian constraints to ensure safe
 83 learning [41, 33]. Our work does not take the Lagrangian approach to learn a safe policy. Instead,
 84 our method explicitly treats the safe policy learning as a sequence modeling problem, similar to the
 85 previous CDT approach [30]. Such an approach enjoys the simplicity regarding the algorithm design,
 86 as well as the robustness to the algorithm performance.

87 3 Preliminaries

88 **Safe Offline RL.** We formulate the environment as a Constrained Markov Decision Process (CMDP),
 89 a mathematical framework for addressing the safe RL problem [2]. A CMDP comprises a tuple
 90 $(\mathcal{S}, \mathcal{A}, \mathcal{P}, r, c, \mu_0)$, where \mathcal{S} denotes the state space, \mathcal{A} signifies the action space, $\mathcal{P} : \mathcal{S} \times \mathcal{A} \times \mathcal{S} \rightarrow$
 91 $[0, 1]$ represents the transition function, $r : \mathcal{S} \times \mathcal{A} \times \mathcal{S} \rightarrow \mathbb{R}$ stands for the reward function, and
 92 $\mu_0 : \mathcal{S} \rightarrow [0, 1]$ indicates the initial state distribution. In CMDP, $c : \mathcal{S} \times \mathcal{A} \times \mathcal{S} \rightarrow [0, C_{\max}]$ quantifies
 93 the cost incurred for violating constraints, where C_{\max} denotes the maximum cost allowable.

94 We denote the policy by $\pi : \mathcal{S} \times \mathcal{A} \rightarrow [0, 1]$, while $\tau = \{s_1, a_1, r_1, c_1, \dots, s_T, a_T, r_T, c_T\}$ delineates
 95 the trajectory containing state, action, and cost information throughout the maximum episode
 96 length T . We use $\tau.s_t, \tau.a_t, \tau.r_t, \tau.c_t$ to denote the t -th state, action, reward and cost in trajectory
 97 τ . For each time step t , the action a_t is drawn following the distribution $\pi(s_t, \cdot)$, and the next
 98 state s_{t+1} is drawn following the transition function $\mathcal{P}(s_t, a_t, \cdot)$. The cumulative reward and cost
 99 for a trajectory τ are represented as $R(\tau) = \sum_{t=1}^T r_t$ and $C(\tau) = \sum_{t=1}^T c_t$. We also denote
 100 $R_t = \sum_{i=t}^T r_i$ by the return-to-go (RTG) at t -th step, $C_t = \sum_{i=t}^T c_i$ by the cost-to-go (CTG)
 101 at t -th step as well. For simplicity, we define $Q_r^\pi(s, a) = \mathbb{E}_{\tau \sim \pi}[R(\tau) | \tau.s_1 = s, \tau.a_1 = a]$ as
 102 the expected return of the policy starting from initial state s and action a . Similarly, we denote
 103 $Q_c^\pi(s, a) = \mathbb{E}_{\tau \sim \pi}[C(\tau) | \tau.s_1 = s, \tau.a_1 = a]$ as the expected cost. The agent’s objective is to
 104 determine a policy π^κ that maximizes the reward while ensuring that the cumulative cost for constraint
 105 violation remains below the threshold κ :

$$\begin{aligned} \pi^\kappa &= \arg \max_{\pi} \mathbb{E}_{\tau \sim \pi, \tau.s_1 \sim \mu_0}[R(\tau)], \\ \text{s.t. } &\mathbb{E}_{\tau \sim \pi, \tau.s_1 \sim \mu_0}[C(\tau)] \leq \kappa. \end{aligned} \quad (1)$$

106 For safe offline RL, the agent learns the optimal safe policy purely from a static dataset \mathcal{T} that is
 107 previously collected with an unknown behavior policy (or policies). Specifically, \mathcal{T} consists of m
 108 episodes τ_i with the maximum length T , which is $\mathcal{T} := \{\tau_1, \dots, \tau_m\}$.

109 **Constrained Decision Transformer.** Our algorithm builds on the Constrained Decision Trans-
 110 former (CDT) [30]. We briefly introduce the details of CDT here, and we leave more details in
 111 the appendix. CDT utilizes the return-conditioned sequential modeling framework to accommo-
 112 date varying constraint thresholds during deployment, ensuring both safety and high return. CDT
 113 employs a stochastic Gaussian policy representation to generate the action at t -th time step, i.e.,
 114 $a_t \sim \pi_\theta(\cdot | o_t) = \mathcal{N}(\mu_\theta(o_t), \Sigma_\theta(o_t))$, where $o_t = \{R_{t-K:t}, C_{t-K:t}, s_{t-K:t}, a_{t-K:t-1}\}$ represents
 115 the truncated history from step $t - K$ to t , $K \in \{1, \dots, t - 1\}$ indicates the context length and
 116 θ denotes the CDT policy parameters. During the training phase, CDT generates the training set
 117 $\{(o_t, a_t)\}$ by splitting trajectories in \mathcal{T} into shorter contexts with length K , then it trains the policy by
 118 minimizing the prediction loss between a_t and $\pi_\theta(\cdot | o_t)$. During the inference phase, CDT selects the
 119 initial return-to-go R_1 as well as the cost-to-go C_1 , then it selects the action a_t based on the current

Algorithm 1 Relabeling trajectory set

Require: Trajectories dataset \mathcal{T} , cost limitation list $\mathcal{K} = [\kappa_1, \kappa_2, \dots, \kappa_m]$, reward critic functions $\mathbb{Q}_r = \emptyset$, cost critic functions $\mathbb{Q}_c = \emptyset$

- 1: //Run CPQ under various cost limitations to derive distinct reward critic and cost critic functions.
- 2: **for** κ_i in \mathcal{K} **do**
- 3: Run CPQ(Algorithm 3) with the constraint limit κ_i , obtain the reward and cost critic Q_r, Q_c
- 4: Set $\mathbb{Q}_r[\kappa_i] = Q_r; \mathbb{Q}_c[\kappa_i] = Q_c$
- 5: **end for**
- 6: //Relabel trajectories in the dataset.
- 7: **for** τ in \mathcal{T} **do**
- 8: Select $\kappa^\tau = \arg \min_{\kappa_i \in \mathcal{K}} |\mathbb{Q}_c[\kappa_i](\tau.s_0, \tau.a_0) - C_0|$
- 9: Set $Q_r^\tau = \mathbb{Q}_r[\kappa^\tau]; Q_c^\tau = \mathbb{Q}_c[\kappa^\tau]$
- 10: Take τ' as the output of Algorithm 2, based on τ, Q_r^τ and $Q_c^\tau, \mathcal{T} \leftarrow \mathcal{T} \cup \{\tau'\}$
- 11: **end for**
- 12: Use Pareto-Frontier augmentation method in Algorithm 4 to augment \mathcal{T}

Ensure: New trajectory dataset \mathcal{T}

120 R_t, C_t . It is worth noting that CDT enables the agent to dynamically adjust its constraint threshold
121 during deployment without using a Lagrangian-type update.

122 **Constraints Penalized Q-learning.** Our algorithm relies on a component that delivers precise value
123 estimates for the state. We take a state-of-the-art method called Constraints Penalized Q-learning
124 (CPQ) [41] as our primary approach and discuss its details as follows. CPQ adopts the actor-critic
125 framework, which maintains four components at each iteration. They are a policy $\pi : \mathcal{S} \times \mathcal{A} \rightarrow [0, 1]$,
126 a discriminator $\nu : \mathcal{S} \times \mathcal{A} \rightarrow \{0, 1\}$ that decides whether a given state-action pair (s, a) are
127 out-of-distribution (OOD) of the behavior policy of the offline dataset; a reward critic function
128 $Q_r : \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}$ and a cost critic function $Q_c : \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}$. During the training phase, CPQ first
129 pre-trains ν by a Conditional Variational Autoencoder (CVAE) [16, 43]. At t -th step, CPQ maintains
130 its current policy π . Then it updates the cost critic first, following

$$Q_c = \arg \min_Q -\alpha \mathbb{E}_{s,a \sim \mathcal{T}} [Q(s, a) \nu(s, a)] \\ + \mathbb{E}_{s,a,s',c \sim \mathcal{T}} [(Q(s, a) - c - \gamma \mathbb{E}_{a' \sim \pi(\cdot|s')} [Q(s', a')])^2],$$

131 where $0 < \gamma < 1, \alpha$ are tunable parameters. Next, CPQ updates the reward critic by optimizing the
132 following cost-penalized Bellman equation, which is

$$Q_r = \arg \min_Q \mathbb{E}_{s,a,s',r \sim \mathcal{T}} [(Q(s, a) - r \\ - \gamma \mathbb{E}_{a' \sim \pi} [\mathbb{I}(Q_c(s', a') \leq \kappa) Q(s', a')])^2], \quad (2)$$

133 where κ is the cost constraint defined in Equation (1). Finally, given Q_r and Q_c , CPQ performs any
134 policy optimization method to obtain π' , which maximizes the following constrained optimization
135 problem:

$$\pi' = \arg \max_{\pi} \mathbb{E}_{s \sim \mathcal{T}} \mathbb{E}_{a \sim \pi(\cdot|s)} [\mathbb{I}(Q_c(s, a) \leq \kappa) Q_r(s, a)]. \quad (3)$$

136 4 Method

137 **Algorithm Overview.** We present a framework called Constrained Q-learning Decision Transformer
138 (CQDT), which exploits the Constrained Dynamic Programming approach, specifically CPQ, to
139 overcome the limitations of CDT. The algorithm details are summarized in Algorithm 1. Intuitively
140 speaking, CQDT takes a trajectory dataset \mathcal{T} as its input and outputs an augmented dataset \mathcal{T}' based
141 on \mathcal{T} . It then trains CDT over the augmented dataset \mathcal{T}' . Next, we describe how to augment \mathcal{T} in
142 steps.

143 **First Step: Training CPQ to Obtain Value Functions.** The objective of this step is to train CPQ to
144 obtain precise estimates of action value functions across different cost limit κ settings. We maintain a
145 list of cost limitations, \mathcal{K} , which includes m different cost limits, $\kappa_1, \dots, \kappa_m$, and run CPQ m times
146 to obtain corresponding Q_r and Q_c functions, denoted as $\mathbb{Q}_r[\kappa_i]$ and $\mathbb{Q}_c[\kappa_i]$, respectively, as outlined

Algorithm 2 Relabeling one trajectory

Require: Trajectory τ , reward critic Q_r^τ and cost critic Q_c^τ
1: Set T as the length of τ , the new trajectory $\tau' = \tau$, $\tau'.R_{T+1} = \tau'.C_{T+1} = 0$
2: **for** $t = T + 1, \dots, 2$ **do**
3: Set $V_r^\tau(\tau.s_t) = Q_r^\tau(\tau.s_t, \tau.a_t)$, $V_c^\tau(\tau.s_t) = Q_c^\tau(\tau.s_t, \tau.a_t)$
4: $\tau'.R_{t-1} \leftarrow \tau.r_{t-1} + \tau'.R_t$, $\tau'.C_{t-1} \leftarrow \tau.c_{t-1} + \tau'.C_t$
5: **if** $V_c^\tau(\tau.s_t) \leq \tau'.C_t$ and $V_r^\tau(\tau.s_t) \geq \tau'.R_t$ **then**
6: $\tau'.R_{t-1} \leftarrow \tau.r_{t-1} + V_r^\tau(\tau.s_t)$, $\tau'.C_{t-1} \leftarrow \tau.c_{t-1} + V_c^\tau(\tau.s_t)$
7: **end if**
8: **end for**
Ensure: Relabeled trajectory τ'

147 in line 4 of Algorithm 1. In the subsequent steps, the \mathbb{Q}_r and \mathbb{Q}_c lists are utilized to relabel the RTG
148 and CTG for each trajectory. For additional details, please refer to Appendix A.1.

149 **Second Step: Relabeling Trajectories.** Now we describe how to relabel a given trajectory τ using
150 \mathbb{Q}_r and \mathbb{Q}_c in detail. The steps are summarized in Algorithm 2. To demonstrate that, recall that our
151 goal is to learn π^κ that maximizes the expected return under the κ constraint. Assume that for the
152 trajectory τ , the policy π that generates τ satisfies the constraint κ . Therefore, in order to further
153 push the agent to learn π^κ instead of only learning π , we generate a new trajectory τ' identical to τ ,
154 with different $\tau'.R_i, \tau'.C_i$, to make τ' similar to a trajectory generated by π^κ . Our strategy is simple:
155 we first replace the last RTG and CTG of τ' as 0. At the t -th step of τ' , we would like to calculate
156 $\tau'.R_{t-1}$ and $\tau'.C_{t-1}$. Then we either to use the existing RTG ($\tau'.R_t$) and CTG ($\tau'.C_t$) to update
157 $\tau'.R_{t-1}$ and $\tau'.C_{t-1}$ (line 4 in Algorithm 2), or we use the learned reward critic and cost critic to
158 update $\tau'.R_{t-1}$ and $\tau'.C_{t-1}$ (line 6 in Algorithm 2), if the reward critic and cost critic provide a
159 more "aggressive" approximation, i.e., $V_c^\tau(\tau.s_t) \leq \tau'.C_t$ and $V_r^\tau(\tau.s_t) \geq \tau'.R_t$ (line 5 in Algorithm
160 2). Here V_r^τ and V_c^τ are learned reward and cost critics, and they are selected from \mathbb{Q}_r and \mathbb{Q}_c
161 to make sure that the cost constraint estimation Q_c^τ is close to the true CTG C_0 (line 8 in Algorithm 2).
162 We summarize the relabeling process in Algorithm 2.

163 The most notable difference between our relabeling strategy and the previous one for offline RL
164 [42] is that we relabel RTG and CTG *jointly and simultaneously*. If we only relabel each trajectory
165 based on their RTG and CTG separately, we might obtain unsafe trajectories, which hurts the overall
166 performance of CQDT. Instead, our strategy ensures that, each safe trajectory will still be safe after
167 relabeling, which is crucial for the safe RL setting. Our experimental results in the later section
168 suggest the effectiveness of our relabeling strategy.

169 **Third Step: Post-Processing Steps for the Final Trajectory.** Now, we have produced an augmented
170 trajectory dataset \mathcal{T} which consists of the original trajectories τ and the new trajectories τ' . Finally,
171 we introduce some additional post-processing steps over \mathcal{T} from existing works [30, 42] for the
172 further performance improvement of CQDT.

173 The first post-processing technique we adopt is to resolve the potential conflict between a high RTG
174 and a low CTG. Due to the nature of safe RL, we would like to first guarantee the safety of our
175 learned policy. Following [30], we use a Pareto Frontier-based data augmentation technique to further
176 generate new trajectories and add them to \mathcal{T} . The second post-processing technique aims to maintain
177 the consistency of the RTG and CTG within the input sequence of CDT. Due to space limitations, we
178 defer the detailed in Appendix A.2.

179 5 Experiment

180 In this section, we begin by outlining the fundamental settings of our experiment. We then show the
181 performance of CQDT under a series of experiments, each addressed a key challenge as follows:

- 182 • How does CQDT compare with CDT and other offline safe reinforcement learning methods in
183 terms of performance? Additionally, how does the choice of the value function affect CQDT's
184 performance?
- 185 • Is CQDT capable of performing effective stitching?

	CarCircle		CarRun		AntRun		DroneCircle		DroneRun		Average	
	reward	cost	reward	cost	reward	cost	reward	cost	reward	cost	reward	cost
BC-Safe*	0.500	0.840	0.940	0.220	0.650	1.090	0.560	0.570	0.280	0.740	0.586	0.692
BCQ-Lag *	0.630	1.890	0.940	0.150	0.760	5.110	0.800	3.070	0.720	5.540	0.770	3.152
BEAR-Lag *	0.740	2.180	0.680	7.780	0.150	0.730	0.780	3.680	0.420	2.470	0.554	3.368
COptiDICE *	0.490	3.140	0.870	0.000	0.610	0.940	0.260	1.020	0.670	4.150	0.580	1.850
CDT*	0.750	0.950	0.990	0.650	0.720	0.910	0.630	0.980	0.630	0.790	0.744	0.856
CPQ *	0.710	0.330	0.950	1.790	0.030	0.020	-0.220	1.280	0.330	3.520	0.360	1.388
CQDT(ours)	0.767 \pm 0.081	0.895 \pm 0.103	0.994 \pm 0.361	0.886 \pm 0.138	0.735 \pm 0.094	0.832 \pm 0.205	0.753 \pm 0.075	0.981 \pm 0.285	0.640 \pm 0.075	0.812 \pm 0.046	0.778 \pm 0.137	0.881 \pm 0.155
BCQL-CDT	0.733 \pm 0.206	0.893 \pm 0.002	0.983 \pm 0.186	1.758 \pm 0.072	0.673 \pm 0.029	0.793 \pm 0.038	0.685 \pm 0.275	0.747 \pm 0.143	0.621 \pm 0.101	0.597 \pm 0.239	0.739 \pm 0.159	0.958 \pm 0.099
BEARL-CDT	0.735 \pm 0.227	0.929 \pm 0.178	0.999 \pm 0.059	1.707 \pm 0.134	0.665 \pm 0.226	0.719 \pm 0.162	0.684 \pm 0.034	0.766 \pm 0.122	0.597 \pm 0.266	0.647 \pm 0.091	0.736 \pm 0.163	0.953 \pm 0.137
Ablation①	0.785 \pm 0.145	0.981 \pm 0.238	0.977 \pm 0.296	1.357 \pm 0.082	0.560 \pm 0.053	0.313 \pm 0.218	0.633 \pm 0.076	0.837 \pm 0.224	0.636 \pm 0.094	0.551 \pm 0.175	0.718 \pm 0.133	0.808 \pm 0.187
Ablation②	0.767 \pm 0.201	0.964 \pm 0.230	0.982 \pm 0.186	1.901 \pm 0.110	0.702 \pm 0.194	0.874 \pm 0.076	0.764 \pm 0.263	1.202 \pm 0.130	0.621 \pm 0.149	0.776 \pm 0.051	0.767 \pm 0.198	1.143 \pm 0.119

Table 1: Evaluation results for normalized reward and cost are provided. **Bold:** Safe agents with a normalized cost smaller than 1. **Gray:** Unsafe agents. **Blue:** Safe agent achieving the highest reward. Each value is averaged over 3 distinct cost thresholds ($\kappa = 10, 20, 40$), 20 episodes and 3 seeds, following the setting in [29]. Results for baselines with * are copied from [29]. BCQL-CDT: Validation of the impact of the value function on algorithm performance: Using the value functions in BCQ-Lag. BEARL-CDT: Validation of the impact of the value function on algorithm performance: Using the value functions in BEAR-Lag. Ablation①: Effect of removing constraints learning by separately considering reward and cost relabeling. Ablation②: Without PF-based augmentation.

- 186 • What is the impact of the various components of CQDT on its overall performance?
- 187 • How does CQDT perform in a sparse reward environment?

188 5.1 Basic Experiment Setting

189 An overview of the basic setup is provided. See Appendix B for details.

190 **Tasks and Environments.** We utilize established safety reinforcement learning tasks within the
191 **BulletSafetyGym** environment following [30]. In our experiment, we focus on two tasks: Circle and
192 Run, and train three types of agents: Car, Ant, and Drone.

193 **Dataset.** We utilize the offline safe RL dataset from [29]. For each environment, the dataset is
194 collected by training different implemented algorithms under gradually varied cost thresholds and
195 collecting the generated trajectories. The algorithms employed during this procedure include CPO
196 [1], FOCOPS [45], CVPO [28], among others.

197 **Baselines and Parameters Setting.** In our selection of baseline models, we have covered a variety
198 of established offline safe RL methods, as thoroughly discussed in [29]. The baselines encompass
199 CDT [30], **Behavior Cloning (BC)** [41], **COptiDICE** [24], **CPQ** [41], **BCQ-Lag** [11, 37], and
200 **BEAR-Lag**[21, 37].

201 5.2 CQDT Performance Analysis

202 **Evaluation Metrics.** We employ the normalized reward return and normalized cost return as our com-
203 parison metrics [10]. Let $R_{\max}(\mathcal{T})$ and $R_{\min}(\mathcal{T})$ denote the maximum and minimum empirical reward
204 returns for task \mathcal{M} , respectively. The normalized reward is calculated as $R_{\text{norm}}(\mathcal{T}) = \frac{R_{\pi}(\mathcal{T}) - R_{\min}(\mathcal{T})}{R_{\max}(\mathcal{T}) - R_{\min}(\mathcal{T})}$,
205 where R_{π} represents the raw return of policy π . For the cost measure, we use the normalized cost
206 return defined as $C_{\text{norm}}(\mathcal{T}) = \frac{C_{\pi}}{\kappa + \epsilon}$, where κ denotes the cost limitation of the task and ϵ is a small
207 positive number incorporated to ensure numerical stability. The term C_{π} represents the evaluated
208 accumulated cost return of policy π . If the cost return surpasses $\kappa + \epsilon$, the normalized cost return
209 exceeds one; otherwise, it remains within the range of one.

210 **Result Analysis.** Table 1 provides a comprehensive overview of the performance of contemporary
211 offline RL strategies in various environments and task configurations. Performance is assessed using
212 reward and cost metrics, with evaluation criteria outlined as follows:

- 213 - When $C_{\text{norm}} \leq 1$, a larger R_{norm} indicates superior strategy performance.
- 214 - When $C_{\text{norm}} \geq 1$, a smaller C_{norm} signifies enhanced strategy performance.
- 215 - Comparing strategies with $C_{\text{norm}} \leq 1$ and $C_{\text{norm}} \geq 1$, preference is given to the policy with
216 $C_{\text{norm}} \leq 1$.

217 Compared to other baselines, our proposed CQDT method achieves maximum return while ensuring
 218 safety. CPQ, BCQ-Lag, and BEAR-Lag, three Q-learning-based safe reinforcement learning methods,
 219 encounter challenges in balancing safety and reward optimization. The BC-Safe method, grounded in
 220 imitation learning and trained on provided data, exhibits suboptimal performance during the test phase.
 221 This phenomenon may be attributed to the scarcity of safe data in the training dataset, indicating a
 222 lack of robustness. COptiDICE employs novel estimators to evaluate policy constraints and achieves
 223 suboptimal rewards, while facing challenges related to adhering to strict safety constraints.

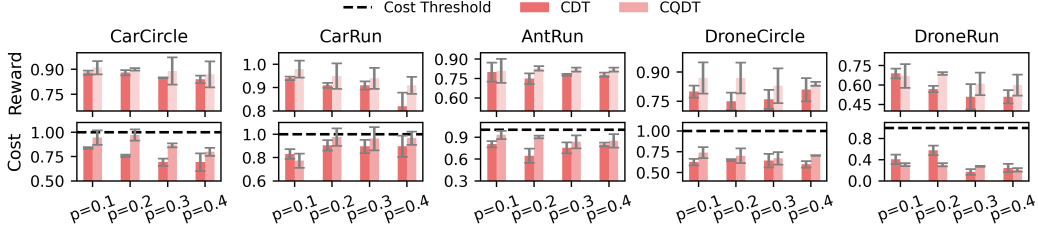


Figure 5: Verification of stitching-reward ability with $p = 0.x$ values representing various suboptimal datasets. Suboptimal datasets were generated by removing safe trajectories that fall within the top $x\%$ of cumulative rewards. Higher p -values indicate the removal of more high-reward safe paths, degrading dataset quality. The first row illustrates cumulative rewards obtained by CQDT and CDT trained on these datasets for different tasks, while the second row shows the corresponding cumulative costs. The black dashed line denotes the cost limitation.

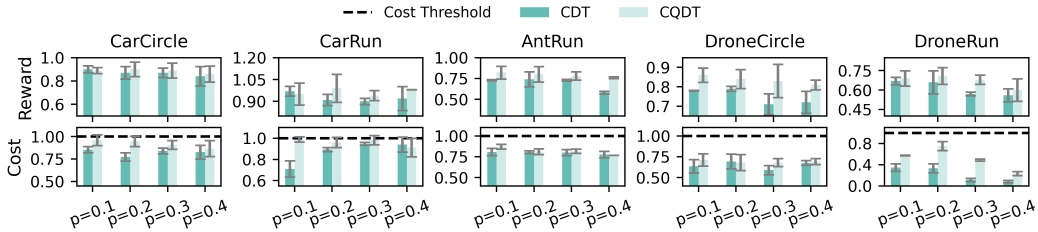


Figure 6: Verification of stitching-cost ability with $p = 0.x$ values corresponding to various suboptimal datasets. These datasets were created by excluding trajectories with low cumulative costs. As the p -value increases, fewer trajectories with small cumulative costs are retained, resulting in increasingly suboptimal datasets.

224 The CQDT method presented in our work leverages additional value functions to relabel trajectories,
 225 enhancing the model’s stitching capability and enabling it to achieve state-of-the-art performance.
 226 To demonstrate the effectiveness of the CQDT method, we conduct a comprehensive comparative
 227 analysis examining the impact of various value functions on algorithmic performance. BCQ-Lag
 228 and BEAR-Lag, both grounded in offline safe reinforcement learning and based on Q-Learning, are
 229 employed for this investigation. The Q_r and Q_c functions in these methodologies are employed to
 230 estimate the RTG and CTG values of the original trajectory. To signify the enhanced versions of
 231 these methods, we specifically refer to them as **BCQL-CDT** and **BEARL-CDT**, respectively. Refer
 232 Appendix C.2 for further details.

233 Our experimental results, detailed in Table 1, indicate that BCQL-CDT and BEARL-CDT perform
 234 similarly to CDT in the selected tasks, although they do not reach the performance of CQDT. The
 235 performance discrepancy between BCQ-Lag and BEAR-Lag compared to CPQ suggests suboptimality
 236 in relabeling the original trajectories using their respective value functions. This contributes to the
 237 varied performance among BCQL-CDT, BEARL-CDT, and CQDT.

238 We also conduct two ablation studies, Ablation ① and Ablation ②, to assess the impact of various
 239 components within CQDT. The results of these experiments are provided in Table 1. For further
 240 details on the ablation studies, please refer to Appendix C.1. In addition, we evaluate the Zero-Shot
 241 Adaptation capability and robustness of CQDT. For more information, refer to Appendix C.4.

242 5.3 The Stitching Capability of CQDT

243 We conduct an evaluation of CQDT’s stitching capabilities for reward and cost by creating various
 244 suboptimal datasets for five tasks and comparing the performance of CQDT and CDT across these
 245 datasets.

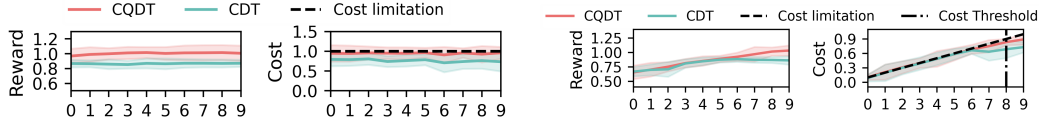


Figure 7: Comparison of CQDT and CDT performance in the Sparse-Reward DroneCircle environment with varying target return and fixed target cost. **Figure 8:** Comparison of CQDT and CDT performance in the Sparse-Reward DroneCircle environment with varying target costs and fixed target return.

246 To evaluate the reward stitching ability, we remove the top $X\%$ of trajectories with the highest RTG
 247 from batches of trajectories. As the value of X increases, more high-return trajectories are excluded
 248 from the dataset. To create a suboptimal dataset for evaluating cost stitching capability, we group
 249 trajectories based on their RTG, then we remove the trajectories with lowest CTG in each group. In
 250 detail, we divide the trajectories into $\frac{Max\ Return}{10}$ groups, where Max Return denotes the highest
 251 return among all trajectories. Within each group, we remove trajectories with the lowest $X\%$ CTG
 252 from each group. Such a setup allows us to detect the stitching ability for a safe offline RL agent, as
 253 we want her to learn safe and high-return policy from unsafe and low-return trajectories. We leave
 254 the detailed parameter setup and the visualization of these suboptimal datasets in the Appendix C.3.

255 We present the performance of CQDT on different suboptimal datasets in Figures 5 and 6. These
 256 experiments demonstrate that as the value of X increases, leading to the removal of higher-quality
 257 trajectories, the performance of policies generated by CQDT and CDT deteriorates. We can observe
 258 that the cumulative reward decreases. However, CQDT consistently outperforms CDT, demonstrating
 259 its superior stitching ability. Even when trained with suboptimal datasets, CQDT effectively utilizes
 260 these datasets to maximize performance by leveraging its stitching capabilities. Superior performance
 261 by CQDT highlights its unique ability to stitch suboptimal trajectories, a capability not present in
 262 CDT. This stitching ability enables CQDT to achieve better overall performance.

263 5.4 Performance of CQDT in Sparse Reward Environment

264 The experiments in the previous sections demonstrate that CQDT performs well in dense reward
 265 environments. In this section, we evaluate and analyze the performance of CQDT in the sparse reward
 266 environment. Since there is no publicly available dataset or corresponding environment for sparse
 267 reward scenarios in the field of offline safe RL, we build our own environment based on existing
 268 datasets. Specifically, we select the existing DroneCircle environment to create a Sparse-Reward
 269 DroneCircle environment.

270 For any trajectory τ in DroneCircle, we aim to create a new trajectory τ' as follows. We consider
 271 each subsequence in τ with length 10, i.e., $\tau.r_{10k}, \tau.r_{10k+1}, \dots, \tau.r_{10k+9}$. We then set $\tau' = \tau$, while
 272 it replaces $\tau'.r_{10k+9}$ with $\tau.r_{10k} + \tau.r_{10k+1} + \dots + \tau.r_{10k+9}$, and replaces $\tau'.r_{10k+i}, 0 \leq i < 9$
 273 with 0. Such an operation keeps the total reward of τ' unchanged, while it makes τ' a trajectory with
 274 sparse rewards. Accordingly, we use the Sparse-Reward DroneCircle Offline Dataset to train CQDT
 275 and CDT. During testing, we adopt the similar strategy, where each agent encounters 0 reward in time
 276 step $10k, 10k + 1, \dots, 10k + 8$, and she encounters $\tau.r_{10k} + \tau.r_{10k+1} + \dots + \tau.r_{10k+9}$ at time step
 277 $10k + 9$. We do not change the cost distribution.

278 Figures 7 and 8 show the performance comparison between CQDT and CDT under different target
 279 return and target cost settings in the Sparse-Reward DroneCircle environment. The results indicate
 280 that CQDT consistently outperforms CDT across various settings of target return and target cost.
 281 These experimental results demonstrate that CQDT maintains its superiority even in sparse-reward
 282 environments.

283 6 Conclusion and Future Work

284 In this work, we proposed the Constrained Q-learning Decision Transformer (CQDT) for safe
 285 offline RL. Our approach replaces reward-to-go and cost-to-go in the training data with dynamic
 286 programming-based learning-based reward return and cost return, which brings the stitching ability
 287 and addresses the weakness of the Constrained Decision Transformer (CDT). Our evaluation shows
 288 that our approach is able to outperform existing safe RL baseline algorithms. One potential future
 289 direction is to build a theoretical analysis to justify the effectiveness of our learning-based constraint
 290 approach to safe RL, similar to previous analyses applied to general goal-based RL algorithms [5].

References

- 291
- 292 [1] Joshua Achiam, David Held, Aviv Tamar, and Pieter Abbeel. Constrained policy optimization.
293 In *International conference on machine learning*, pages 22–31. PMLR, 2017.
- 294 [2] Eitan Altman. Constrained markov decision processes with total cost criteria: Lagrangian
295 approach and dual linear program. *Mathematical methods of operations research*, 48:387–417,
296 1998.
- 297 [3] Anurag Arnab, Mostafa Dehghani, Georg Heigold, Chen Sun, Mario Lučić, and Cordelia
298 Schmid. Vivit: A video vision transformer. In *Proceedings of the IEEE/CVF international
299 conference on computer vision*, pages 6836–6846, 2021.
- 300 [4] Iz Beltagy, Matthew E Peters, and Arman Cohan. Longformer: The long-document transformer.
301 *arXiv preprint arXiv:2004.05150*, 2020.
- 302 [5] David Brandfonbrener, Alberto Bietti, Jacob Buckman, Romain Laroche, and Joan Bruna. When
303 does return-conditioned supervised learning work for offline reinforcement learning? *Advances
304 in Neural Information Processing Systems*, 35:1542–1553, 2022.
- 305 [6] Greg Brockman, Vicki Cheung, Ludwig Pettersson, Jonas Schneider, John Schulman, Jie Tang,
306 and Wojciech Zaremba. Openai gym. *arXiv preprint arXiv:1606.01540*, 2016.
- 307 [7] Lili Chen, Kevin Lu, Aravind Rajeswaran, Kimin Lee, Aditya Grover, Misha Laskin, Pieter
308 Abbeel, Aravind Srinivas, and Igor Mordatch. Decision transformer: Reinforcement learning
309 via sequence modeling. *Advances in neural information processing systems*, 34:15084–15097,
310 2021.
- 311 [8] Yinlam Chow, Ofir Nachum, Aleksandra Faust, Edgar Duenez-Guzman, and Mohammad
312 Ghavamzadeh. Lyapunov-based safe policy optimization for continuous control. *arXiv preprint
313 arXiv:1901.10031*, 2019.
- 314 [9] Yan Duan, Xi Chen, Rein Houthoofd, John Schulman, and Pieter Abbeel. Benchmarking deep
315 reinforcement learning for continuous control. In *International conference on machine learning*,
316 pages 1329–1338. PMLR, 2016.
- 317 [10] Justin Fu, Aviral Kumar, Ofir Nachum, George Tucker, and Sergey Levine. D4rl: Datasets for
318 deep data-driven reinforcement learning. *arXiv preprint arXiv:2004.07219*, 2020.
- 319 [11] Scott Fujimoto, David Meger, and Doina Precup. Off-policy deep reinforcement learning
320 without exploration. In *International conference on machine learning*, pages 2052–2062.
321 PMLR, 2019.
- 322 [12] Javier Garcia and Fernando Fernández. A comprehensive survey on safe reinforcement learning.
323 *Journal of Machine Learning Research*, 16(1):1437–1480, 2015.
- 324 [13] Omer Gottesman, Fredrik Johansson, Joshua Meier, Jack Dent, Donghun Lee, Srivatsan Srinivasan,
325 Linying Zhang, Yi Ding, David Wihl, Xuefeng Peng, et al. Evaluating reinforcement
326 learning algorithms in observational health settings. *arXiv preprint arXiv:1805.12298*, 2018.
- 327 [14] Shangding Gu, Long Yang, Yali Du, Guang Chen, Florian Walter, Jun Wang, Yaodong Yang,
328 and Alois Knoll. A review of safe reinforcement learning: Methods, theory and applications.
329 *arXiv preprint arXiv:2205.10330*, 2022.
- 330 [15] Tuomas Haarnoja, Aurick Zhou, Pieter Abbeel, and Sergey Levine. Soft actor-critic: Off-
331 policy maximum entropy deep reinforcement learning with a stochastic actor. In *International
332 conference on machine learning*, pages 1861–1870. PMLR, 2018.
- 333 [16] Boris Ivanovic, Karen Leung, Edward Schmerling, and Marco Pavone. Multimodal deep
334 generative models for trajectory prediction: A conditional variational autoencoder approach.
335 *IEEE Robotics and Automation Letters*, 6(2):295–302, 2020.
- 336 [17] Michael Janner, Justin Fu, Marvin Zhang, and Sergey Levine. When to trust your model:
337 Model-based policy optimization. *Advances in neural information processing systems*, 32,
338 2019.

- 339 [18] Michael Janner, Qiyang Li, and Sergey Levine. Offline reinforcement learning as one big
340 sequence modeling problem. *Advances in neural information processing systems*, 34:1273–
341 1286, 2021.
- 342 [19] Dmitry Kalashnikov, Alex Irpan, Peter Pastor, Julian Ibarz, Alexander Herzog, Eric Jang,
343 Deirdre Quillen, Ethan Holly, Mrinal Kalakrishnan, Vincent Vanhoucke, et al. Scalable deep
344 reinforcement learning for vision-based robotic manipulation. In *Conference on robot learning*,
345 pages 651–673. PMLR, 2018.
- 346 [20] Rahul Kidambi, Aravind Rajeswaran, Praneeth Netrapalli, and Thorsten Joachims. Morel:
347 Model-based offline reinforcement learning. *Advances in neural information processing systems*,
348 33:21810–21823, 2020.
- 349 [21] Aviral Kumar, Justin Fu, Matthew Soh, George Tucker, and Sergey Levine. Stabilizing off-
350 policy q-learning via bootstrapping error reduction. *Advances in Neural Information Processing*
351 *Systems*, 32, 2019.
- 352 [22] Aviral Kumar, Aurick Zhou, George Tucker, and Sergey Levine. Conservative q-learning for
353 offline reinforcement learning. *Advances in Neural Information Processing Systems*, 33:1179–
354 1191, 2020.
- 355 [23] Hoang Le, Cameron Voloshin, and Yisong Yue. Batch policy learning under constraints. In
356 *International Conference on Machine Learning*, pages 3703–3712. PMLR, 2019.
- 357 [24] Jongmin Lee, Cosmin Paduraru, Daniel J Mankowitz, Nicolas Heess, Doina Precup, Kee-Eung
358 Kim, and Arthur Guez. Coptidice: Offline constrained reinforcement learning via stationary
359 distribution correction estimation. *arXiv preprint arXiv:2204.08957*, 2022.
- 360 [25] Sergey Levine, Aviral Kumar, George Tucker, and Justin Fu. Offline reinforcement learning:
361 Tutorial, review, and perspectives on open problems. *arXiv preprint arXiv:2005.01643*, 2020.
- 362 [26] Timothy P Lillicrap, Jonathan J Hunt, Alexander Pritzel, Nicolas Heess, Tom Erez, Yuval Tassa,
363 David Silver, and Daan Wierstra. Continuous control with deep reinforcement learning. *arXiv*
364 *preprint arXiv:1509.02971*, 2015.
- 365 [27] Ze Liu, Yutong Lin, Yue Cao, Han Hu, Yixuan Wei, Zheng Zhang, Stephen Lin, and Baining
366 Guo. Swin transformer: Hierarchical vision transformer using shifted windows. In *Proceedings*
367 *of the IEEE/CVF international conference on computer vision*, pages 10012–10022, 2021.
- 368 [28] Zuxin Liu, Zhepeng Cen, Vladislav Isenbaev, Wei Liu, Steven Wu, Bo Li, and Ding Zhao.
369 Constrained variational policy optimization for safe reinforcement learning. In *International*
370 *Conference on Machine Learning*, pages 13644–13668. PMLR, 2022.
- 371 [29] Zuxin Liu, Zijian Guo, Haohong Lin, Yihang Yao, Jiacheng Zhu, Zhepeng Cen, Hanjiang
372 Hu, Wenhao Yu, Tingnan Zhang, Jie Tan, et al. Datasets and benchmarks for offline safe
373 reinforcement learning. *arXiv preprint arXiv:2306.09303*, 2023.
- 374 [30] Zuxin Liu, Zijian Guo, Yihang Yao, Zhepeng Cen, Wenhao Yu, Tingnan Zhang, and Ding
375 Zhao. Constrained decision transformer for offline safe reinforcement learning. *arXiv preprint*
376 *arXiv:2302.07351*, 2023.
- 377 [31] Martin Pecka and Tomas Svoboda. Safe exploration techniques for reinforcement learning—an
378 overview. In *Modelling and Simulation for Autonomous Systems: First International Workshop,*
379 *MESAS 2014, Rome, Italy, May 5-6, 2014, Revised Selected Papers 1*, pages 357–375. Springer,
380 2014.
- 381 [32] Xue Bin Peng, Aviral Kumar, Grace Zhang, and Sergey Levine. Advantage-weighted regression:
382 Simple and scalable off-policy reinforcement learning. *arXiv preprint arXiv:1910.00177*, 2019.
- 383 [33] Nicholas Polosky, Bruno C Da Silva, Madalina Fiterau, and Jithin Jagannath. Constrained offline
384 policy optimization. In *International Conference on Machine Learning*, pages 17801–17810.
385 PMLR, 2022.

- 386 [34] Alex Ray, Joshua Achiam, and Dario Amodei. Benchmarking safe exploration in deep rein-
387 forcement learning. *arXiv preprint arXiv:1910.01708*, 7(1):2, 2019.
- 388 [35] John Schulman, Sergey Levine, Pieter Abbeel, Michael Jordan, and Philipp Moritz. Trust
389 region policy optimization. In *International conference on machine learning*, pages 1889–1897.
390 PMLR, 2015.
- 391 [36] Shai Shalev-Shwartz, Shaked Shammah, and Amnon Shashua. Safe, multi-agent, reinforcement
392 learning for autonomous driving. *arXiv preprint arXiv:1610.03295*, 2016.
- 393 [37] Adam Stooke, Joshua Achiam, and Pieter Abbeel. Responsive safety in reinforcement learning
394 by pid lagrangian methods. In *International Conference on Machine Learning*, pages 9133–9143.
395 PMLR, 2020.
- 396 [38] Ashish Vaswani. Attention is all you need. *arXiv preprint arXiv:1706.03762*, 2017.
- 397 [39] Qing Wang, Jiechao Xiong, Lei Han, Han Liu, Tong Zhang, et al. Exponentially weighted
398 imitation learning for batched historical data. *Advances in Neural Information Processing*
399 *Systems*, 31, 2018.
- 400 [40] Thomas Wolf, Lysandre Debut, Victor Sanh, Julien Chaumond, Clement Delangue, Anthony
401 Moi, Pierric Cistac, Tim Rault, Rémi Louf, Morgan Funtowicz, et al. Transformers: State-of-the-
402 art natural language processing. In *Proceedings of the 2020 conference on empirical methods in*
403 *natural language processing: system demonstrations*, pages 38–45, 2020.
- 404 [41] Haoran Xu, Xianyuan Zhan, and Xiangyu Zhu. Constraints penalized q-learning for safe offline
405 reinforcement learning. In *Proceedings of the AAAI Conference on Artificial Intelligence*,
406 volume 36, pages 8753–8760, 2022.
- 407 [42] Taku Yamagata, Ahmed Khalil, and Raul Santos-Rodriguez. Q-learning decision transformer:
408 Leveraging dynamic programming for conditional sequence modelling in offline rl. In *Internat-*
409 *ional Conference on Machine Learning*, pages 38989–39007. PMLR, 2023.
- 410 [43] Jingkang Yang, Kaiyang Zhou, Yixuan Li, and Ziwei Liu. Generalized out-of-distribution
411 detection: A survey. *arXiv preprint arXiv:2110.11334*, 2021.
- 412 [44] Tianhe Yu, Aviral Kumar, Rafael Rafailov, Aravind Rajeswaran, Sergey Levine, and Chelsea
413 Finn. Combo: Conservative offline model-based policy optimization. *Advances in neural*
414 *information processing systems*, 34:28954–28967, 2021.
- 415 [45] Yiming Zhang, Quan Vuong, and Keith Ross. First order constrained optimization in policy
416 space. *Advances in Neural Information Processing Systems*, 33:15338–15349, 2020.

417 **A Implementation Details**

418 In this section, we introduce more details of the implementation of CQDT.

419 **A.1 Implementation Details of CPQ (First Step of CQDT)**

420 **OOD discriminator** ν . The process of learning the OOD action generation distribution ν uses an
 421 OOD detection method based on the Conditional Variational Autoencoder (CVAE). In detail, ν is
 422 based on a decoder $p : \mathcal{S} \times \mathcal{Z} \times \mathcal{A} \rightarrow [0, 1]$ that generates the action following the distribution
 423 $p(s, z, \cdot)$, where $z \in \mathcal{Z} \subseteq \mathbb{R}$ is the hidden state; an encoder $q : \mathcal{S} \times \mathcal{A} \times \mathcal{Z} \rightarrow [0, 1]$ that generates the
 424 latent state following the distribution $q(s, a, \cdot)$. p and q are trained by solving the following evidence
 425 lower bound (ELBO) objective, which is

$$\max_{p,q} \mathbb{E}_{s,a \sim \mathcal{T}} \left[\mathbb{E}_{z \sim q} \log p(s, z, a) - \beta \text{KL}(q(s, a, \cdot) \| N(0, 1)) \right], \quad (4)$$

426 where KL denotes the KL divergence and β is the penalty parameter. We then set ν as

$$\nu(s, a) = \begin{cases} 1 & \text{KL}(q(s, a, \cdot) \| N(0, 1)) \geq d \\ 0 & \text{Otherwise} \end{cases} \quad (5)$$

427 Here we provide a detailed implementation code for CPQ, following (author?) [41].

Algorithm 3 CPQ

Require: Trajectories dataset \mathcal{T} ; constraint limitation κ ; initialize encoder q and decoder p ; training steps M and N .

1: // VAE training

2: **for** $t = 0$ to M **do**

3: Sample mini-batch of state-action pairs $(s, a) \sim \mathcal{T}$ and update p, q through (4)

4: **end for**

5: Update ν following (5)

6: // Policy training

7: Initialize reward critic Q_r , cost critic Q_c , actor π_θ

8: **for** $t = 0$ to N **do**

9: Update cost critic by

$$Q_c = \arg \min_Q -\alpha \mathbb{E}_{s,a \sim \mathcal{T}} [Q(s, a) \nu(s, a)] + \mathbb{E}_{s,a,s',c \sim \mathcal{T}} [(Q(s, a) - c - \gamma \mathbb{E}_{a' \sim \pi(\cdot|s')} [Q(s', a')])^2],$$

10: Update reward critic by

$$Q_r = \arg \min_Q \mathbb{E}_{s,a,s',r \sim \mathcal{T}} [(Q(s, a) - r - \gamma \mathbb{E}_{a' \sim \pi} [\mathbb{I}(Q_c(s', a') \leq \kappa) Q(s', a')])^2],$$

11: Update policy as

$$\pi' = \arg \max_{\pi} \mathbb{E}_{s \sim \mathcal{T}} \mathbb{E}_{a \sim \pi(\cdot|s)} [\mathbb{I}(Q_c(s, a) \leq \kappa) Q_r(s, a)].$$

12: **end for**

Ensure: reward critic $Q_r(s, a)$, cost critic $Q_c(s, a)$

428 **A.2 Implementation details of Data Augmentation (Third Step of CQDT)**

429 We adopt a data augmentation technique based on the Pareto Frontier (PF) from [30]. The dataset
 430 used in safe offline RL is characterized by the PF value $\text{PF}(\kappa, \mathcal{T})$, which is

$$\text{PF}(\kappa, \mathcal{T}) = \max_{\tau \in \mathcal{T}} R(\tau), \quad \text{s.t. } C(\tau) \leq \kappa$$

431 Given that CQDT shares the fundamental framework with CDT, it adopts the target returns-
 432 conditioned policy structure. Consequently, the agent’s behavior becomes sensitive to choices re-
 433 garding target return and cost. This sensitivity limits the valid choices for target cost and reward

Table 2: The hyperparameter configurations for dataset collection algorithms across diverse tasks. The parameters *Cost Start*, *Cost End*, *Epoch Start*, and *Epoch End* are utilized to modify the cost limitations.

Task	Cost Start	Cost End	Epoch Start	Epoch End	Epoch Number	Max Trajectory Length
CarCircle	5	100	100	900	1000	1500
CarRun	5	100	100	900	1000	1500
AntRun	5	200	400	2500	2600	2000
DroneCircle	5	150	500	2500	2600	2000
DroneRun	100	5	50	800	1000	1500

434 return pairs to the PF points. To address this limitation, CQDT employs the same data enhancement
 435 strategy as CDT to improve the relabeled data. We list the details of the data augmentation technique
 436 in Algorithm 4.

Algorithm 4 PF augmentation

Require: Trajectories dataset \mathcal{T} , iteration number N , max length T of trajectories in \mathcal{T}

- 1: Set $c_{\min} = \min_{\tau \in \mathcal{T}} C(\tau)$, $c_{\max} = \max_{\tau \in \mathcal{T}} C(\tau)$, $r_{\max} = \max_{\tau \in \mathcal{T}} R(\tau)$
- 2: **for** $i = 1, \dots, N$ **do**
- 3: $\kappa_i \sim \text{Uniform}(c_{\min}, c_{\max})$ // Sample a cost return
- 4: $\rho_i \sim \text{Uniform}(\text{PF}(\kappa_i, \mathcal{T}), r_{\max})$ // Sample a reward return above the PF value
- 5: $\tau_i^* \leftarrow \arg \max_{\tau \in \mathcal{T}} R(\tau)$, s.t. $C(\tau) \leq \kappa_i$ // Find the closest and safe Pareto trajectory
- 6: Generate $\hat{\tau}_i$, where $\hat{\tau}_i.R_t \leftarrow \tau_i^*.R_t + \rho_i - R(\tau_i^*)$, $\hat{\tau}_i.C_t \leftarrow \tau_i^*.C_t + \kappa_i - C(\tau_i^*)$ for all $1 \leq t \leq T$ // Relabel the reward and cost return
- 7: $\mathcal{T} \leftarrow \mathcal{T} \cup \{\hat{\tau}_i\}$ // Append the trajectory to the dataset
- 8: **end for**

Ensure: Augmented trajectories dataset \mathcal{T}

437 Finally, we generate the final sliced trajectories dataset \mathcal{T}^K which includes consistent trajectory τ
 438 with length K . To generate it, for each $\tau \in \mathcal{T}$, we regenerate new trajectories dataset $\tau^1, \dots, \tau^{T-K}$ as
 439 follows. For τ^t , we set its last RTG and CTG as $\tau^t.R_{t+K} \leftarrow \tau.R_{t+K}$ and $\tau^t.C_{t+K} \leftarrow \tau.C_{t+K}$. Then
 440 we repeatedly apply $\tau^t.R_i \leftarrow \tau.R_i + \tau^t.R_{i+1}$ and $\tau^t.C_i \leftarrow \tau.C_i + \tau^t.C_{i+1}$ for $i = t + K - 1, \dots, t$.
 441 We summarize it in Algorithm 5.

442 B Experiment Setting and Hyperparameters

443 B.1 Tasks Description

444 In our experiments, we employ the BulletSafetyGym environment, which is a suite built atop the
 445 PyBullet physics simulator, resembles *SafetyGym* but features shorter horizons and a larger number
 446 of agents [1, 34, 9, 8, 6]. Within BulletSafetyGym, we deploy three distinct agent models: the Car,
 447 Ant, and Drone. The Ant agent is designed as a four-legged creature with a spherical torso, while the
 448 Car agent, inspired by MIT’s Racecar, features a four-wheeled configuration. The Drone agent is an
 449 aerial vehicle based on the AscTec Hummingbird quadrotor. These agents are employed to complete
 450 the Circle and Run tasks. In the Circle task, they move clockwise on a circle. The reward is defined as

$$r(s) = \frac{v^T[-y, x]}{1 + 3|r_{\text{agent}} - r_{\text{circle}}|} \quad (6)$$

451 (6) is maximized when agents move quickly in a clockwise direction. Costs are incurred if the agent
 452 leaves the safety zone defined by the two yellow boundaries, i.e., $c(s) = \mathbb{I}[|x| > x_{\text{lim}}]$, where x is
 453 the position on the x-axis. In the Run task, agents earn rewards for navigating an avenue located
 454 between two non-physical, penetrable safety boundaries that incur costs upon breach. Additional
 455 costs are applied if agents surpass an agent-specific velocity threshold.

456 B.2 Dataset Collection

457 The dataset collection process is the same as [29], and we include it for completeness. In the data
 458 collection process, we use various algorithms and distinct cost thresholds tailored to each envi-

Algorithm 5 Consistent RTG and CTG

Require: Trajectories dataset \mathcal{T} , context length K

```
1: Set  $\mathcal{T}^K = \emptyset$ 
2: for  $\tau \in \mathcal{T}$  do
3:   for  $t = 1, \dots, T - K$  do
4:     Set  $\tau^t = \{\tau.s_t, \tau.a_t, \tau.r_t, \tau.c_t, \dots, \tau.s_{t+K}, \tau.a_{t+K}, \tau.r_{t+K}, \tau.c_{t+K}\}$ 
5:     Set  $\tau^t.R_{t+K} \leftarrow \tau.R_{t+K}$  and  $\tau^t.C_{t+K} \leftarrow \tau.C_{t+K}$ 
6:     for  $i = t + K - 1, \dots, t$  do
7:       Set  $\tau^t.R_i \leftarrow \tau.r_i + \tau^t.R_{i+1}$  and  $\tau^t.C_i \leftarrow \tau.c_i + \tau^t.C_{i+1}$ 
8:     end for
9:      $\mathcal{T}^K \leftarrow \mathcal{T}^K \cup \{\tau^t\}$ 
10:  end for
11: end for
Ensure: Sliced trajectories dataset  $\mathcal{T}^K$ 
```

Table 3: Description of the experimental datasets: *Max TS* denotes the maximum length of an episode, while *Act. Space* and *State Space* represent the dimensions of action and state, respectively. *MaxCost*, *MinCost*, *MaxReward*, and *MinReward* represent the cumulative reward and cumulative cost obtained by each trajectory. *Traj.* indicates the number of trajectories in the dataset.

Bench.	Task	Max TS	Act. Space	State Space	MaxCost	MinCost	MaxReward	MinReward	Traj.
BulletSafetyGym	CarCircle	300	2	8	100	0	534.306	3.484	1450
	CarRun	200	2	7	40	0	574.653	204.287	651
	AntRun	200	8	33	150	0	955.481	0	1816
	DroneCircle	300	4	18	100	0	996.389	207.79	1923
	DroneRun	200	4	17	140	0	682.83	10.557	1990

459 ronment. The algorithms utilized for dataset acquisition include CPO, FOCOPS, PPOLagrangian,
460 TRPOLagrangian, DDPGLagrangian, SACLagrangian, and CVPO. Notably, PPOLagrangian, TR-
461 POLagrangian, DDPGLagrangian, and SACLagrangian are composite methods combining PPO [1],
462 TRPO [35], DDPG [26], SAC [15], and PID Lagrangian [37] techniques, respectively. Among these
463 algorithms, the first four belong to the category of On-Policy algorithms, while DDPGLagrangian
464 and SACLagrangian fall under Off-On-Policy algorithms. On the other hand, CVPO is classified as
465 an Off-Policy algorithm. Table 2 outlines the hyperparameters utilized for collecting datasets across
466 different tasks, and Table 3 presents detailed information about the constructed datasets.

467 B.3 Hyperparameters

468 We list the hyperparameters for CQDT here as well as the hyperparameters employed in the training
469 of CPQ. Specifically, within CPQ, we utilize the Q_r and Q_c functions for trajectory relabeling. The
470 hyperparameters for CPQ training are listed in Table 4. Meanwhile, the hyperparameters for CDT
471 and CQDT training are comprehensively detailed in Table 5.

472 C Result Details and Discussions

473 C.1 Ablation Study

474 In assessing the impact of various components within CQDT, we conducted two ablation studies.

475 **Joint Relabeling.** This is denoted as Ablation ①, aiming to study whether our adopted relabeling
476 strategy in Algorithm 2 is necessary. We compare it with an alternative relabeling approach which
477 does not take both the reward and cost into consideration together. Instead, the alternative relabeling
478 approach relabels the RTG and CTG similar to what QDT does, which relabels them separately. For
479 a trajectory τ , the alternative approach obtains V_r^τ and V_c^τ the same as Algorithm 2, then for each
480 $t = T + 1, \dots, 2$, it relabels $\tau.R_{t-1} \rightarrow \tau.r_{t-1} + \max(\tau.R_t, V_r^\tau(\tau.s_t))$ when $V_r^\tau(s_{\mathcal{L}}) \geq \tau.R_{\mathcal{L}}$, and it
481 relabels $\tau.C_{t-1} \rightarrow \tau.c_{t-1} + \min(\tau.C_t, V_c^\tau(\tau.s_t))$ when $V_c^\tau(s_t) \leq \tau.C_t$. The main difference between
482 such a relabeling strategy and our adopted strategy for CQDT is the separate consideration for reward
483 and cost relabeling. Our result reveals that with such a relabeling process, CQDT performance is
484 notably worse in scenarios such as AntRun and DroneCircle, which exhibit poor CPQ performance

Table 4: Parameters utilized in training CPQ. $range(x : y, z)$ means the cost limitation begins at x and increases by z each step until it approaches y .

	CarCircle	CarRun	AntRun	DroneCircle	DroneRun
<i>state_dim</i>	8	7	33	18	17
<i>action_dim</i>	2	2	8	4	4
<i>max_action</i>	1.0	1.0	1.0	1.0	1.0
<i>actor_hidden_size</i>	[256, 256]	[256, 256]	[256, 256]	[256, 256]	[256, 256]
<i>critic_hidden_size</i>	[256, 256]	[256, 256]	[256, 256]	[256, 256]	[256, 256]
<i>VAE_hidden_size</i>	400	400	400	400	400
<i>sample_action_num</i>	10	10	10	10	10
<i>Q_r number</i>	2	2	2	2	2
<i>Q_c number</i>	2	2	2	2	2
<i>episode_len</i>	300	200	200	300	200
<i>batchsize</i>	2048	2048	2048	2048	2048
<i>update_steps</i>	100000	100000	100000	100000	100000
<i>vae_lr</i>	0.001	0.001	0.001	0.001	0.001
<i>critic_lr</i>	0.001	0.001	0.001	0.001	0.001
<i>actor_lr</i>	0.0001	0.0001	0.0001	0.0001	0.0001
<i>cost_limitation</i>	(10 : 100, 10)	(5 : 45, 5)	(15 : 150, 15)	(10 : 110, 10)	(15 : 135, 15)

Table 5: Parameters utilized in training CDT and CQDT. Shared parameters are denoted in black fonts, while the independent parameters of CQDT are marked in blue fonts.

	CarCircle	CarRun	AntRun	DroneCircle	DroneRun
<i>state_dim</i>	8	7	33	18	17
<i>action_dim</i>	2	2	8	4	4
<i>max_action</i>	1.0	1.0	1.0	1.0	1.0
<i>embedding_dim</i>	128	128	128	128	128
<i>seq_len</i>	10	10	10	10	10
<i>episode_len</i>	300	200	200	300	200
<i>num_layers</i>	3	3	3	3	3
<i>num_heads</i>	8	8	8	8	8
<i>attention_dropout</i>	0.1	0.1	0.1	0.1	0.1
<i>residual_dropout</i>	0.1	0.1	0.1	0.1	0.1
<i>embedding_dropout</i>	0.1	0.1	0.1	0.1	0.1
<i>action_head_layers</i>	1	1	1	1	1
<i>target&cost_return</i>	[500, (0 : 120, 10)]	[575, (0 : 50, 5)]	[900, (0 : 200, 20)]	[900, (20 : 120, 20)]	[600, (0 : 200, 20)]
<i>target&cost_return</i>	[(0 : 600, 50), 90]	[(70 : 700, 70), 35]	[(0 : 1100, 100), 130]	[(300 : 1400, 100), 90]	[(0 : 1000, 100), 120]
<i>batchsize</i>	2048	2048	2048	2048	2048
<i>learning_rate</i>	0.0001	0.0001	0.0001	0.0001	0.0001
<i>update_steps</i>	100000	100000	100000	100000	100000
<i>weight_decay</i>	0.0001	0.0001	0.0001	0.0001	0.0001

485 when cost return and reward return are considered independently. This indicates that the accuracy of
486 value function predictions significantly influences model performance during the relabeling process,
487 while our proposed CQDT method effectively minimizes the negative impacts caused by inaccuracies
488 in value function predictions.

489 **PF Augmentation.** The second study, denoted as Ablation ②, evaluates CQDT without the PF-based
490 augmentation step [30]. Our results reveal that the PF-based augmentation is essential for ensuring
491 the model’s security under various target cost settings. For instance, in the CarRun task, the absence
492 of the PF augmentation technique results in a policy that exceeds cost limitations.

493 C.2 Detailed Results from Value Function Experiment

494 In this section, we explore the utilization of diverse value functions for state value estimation,
495 substituting RTG and CTG in the original trajectories. We then retrain our CQDT model with the
496 replaced dataset and analyze the differences in results. We systematically investigate the incorporation
497 of Q_r and Q_c as value functions for state value estimation in CPQ, BCQ-Lagrangian (BCQL), and
498 BEAR-Lagrangian (BEARL). The algorithmic details of CPQ are outlined in Algorithm 3. Combining
499 BCQ and BEAR with the Lagrangian approach, which utilizes adaptive penalty coefficients to enforce
500 constraints, results in the formulation of BCQ-Lagrangian and BEAR-Lagrangian.

501 **Implementation Details of BCQ** For details regarding the BCQ, please refer to Algorithm 6. BCQ
502 requires a tuple dataset $B = \{(s, a, r, s')\}$ and returns policy π and reward critic Q . Here, π is

503 defined based on n sampled action a_i , where

$$\pi(s) = \arg \max_{a_i + \xi_\phi(s, a_i, \Phi)} Q_{\theta_1}(s, a_i + \xi_\phi(s, a_i, \Phi)), \{a_i \sim G_w(s)\}_{i=1}^n. \quad (7)$$

Algorithm 6 BCQ Training

Require: Dataset B , iteration number M , target network update rate τ , mini-batch size N , max perturbation Φ , number of sampled actions n , minimum weighting λ .

- 1: Initialize Q -networks $Q_{\theta_1}, Q_{\theta_2}$, perturbation network ξ_ϕ , and VAE $G_w = \{E_{\omega_1}, D_{\omega_2}\}$, with random parameters $\theta_1, \theta_2, \phi, \omega$, and target networks $Q'_{\theta_1}, Q'_{\theta_2}, \xi'_\phi$ with $\theta'_1 \leftarrow \theta_1, \theta'_2 \leftarrow \theta_2, \phi' \leftarrow \phi$.
 - 2: **for** $t = 1, \dots, M$ **do**
 - 3: Sample mini-batch of N transitions (s, a, r, s') from B
 - 4: $\mu, \sigma = E_{\omega_1}(s, a), \tilde{a} = D_{\omega_2}(s, z), z \sim \mathcal{N}(\mu, \sigma)$
 - 5: $\omega \leftarrow \arg \min_\omega \sum (a - \tilde{a})^2 + \text{KL}(\mathcal{N}(\mu, \sigma) \parallel \mathcal{N}(0, 1))$
 - 6: Sample n actions: $\{\tilde{a}_i\}_{i=1}^n \sim G_w(s')$
 - 7: Perturb each action: $\{\tilde{a}_i = a_i + \xi_\phi(s', a_i, \Phi)\}_{i=1}^n$
 - 8: Set value target $y = r + \max_i [\lambda \min_{j=1,2} Q_{\theta'_j}(s', \tilde{a}_i) + (1 - \lambda) \max_{j=1,2} Q_{\theta_j}(s', \tilde{a}_i)]$
 - 9: $\theta \leftarrow \arg \min_\theta \sum (y - Q_\theta(s, a))^2$ and $\phi \leftarrow \arg \max_\phi \sum Q_{\theta_1}(s, a + \xi_\phi(s, a, \Phi))$
 - 10: Update target networks: $\theta'_i \leftarrow \tau \theta_i + (1 - \tau) \theta'_i, \phi' \leftarrow \tau \phi + (1 - \tau) \phi'$
 - 11: **end for**
- Ensure:** Policy π follows (7), reward critic $Q = Q_{\theta_1}$
-

504 **Implementation Details of BEAR** For BEAR algorithm, we put its details in Algorithm 7. BEAR
 505 takes a tuple dataset $B = \{(s, a, r, s')\}$ as its input. It operates based on the maximum mean
 506 discrepancy (MMD) distance as follows

$$\text{MMD}^2(\{x_1, \dots, x_n\}, \{y_1, \dots, y_m\}) = \frac{1}{n^2} \sum_{i, i'} k(x_i, x_{i'}) - \frac{2}{nm} \sum_{i, j} k(x_i, y_j) + \frac{1}{m^2} \sum_{j, j'} k(y_j, y_{j'}). \quad (8)$$

Algorithm 7 BEAR Training

Require: Dataset B , target network update rate τ , iteration number M , sampled actions for MMD n , ratio parameter λ

- 1: Initialize Q -ensemble $\{Q_{\theta_i}\}_{i=1}^K$, actor π_ϕ , Lagrange multiplier α , target networks $\{Q'_{\theta_i}\}_{i=1}^K$, and a target actor π'_ϕ with $\phi' \leftarrow \phi, \theta'_i \leftarrow \theta_i$
- 2: **for** $t = 1, \dots, M$ **do**
- 3: Sample mini-batch of transitions $(s, a, r, s') \sim B$
- 4: Sample p action samples, $\{a_i \sim \pi'_\phi(s')\}_i^p$
- 5: Define $y(s, a) = \max_{a_i} [\min_{j=1, \dots, K} Q'_{\theta_j}(s', a_i) + (1 - \lambda) \max_{j=1, \dots, K} Q_{\theta_j}(s', a_i)]$
- 6: $\forall i, \theta_i \leftarrow \arg \min_{\theta_i} ((Q_{\theta_i}(s, a) - (r + \gamma y(s, a)))^2)$
- 7: Sample actions $\{a_i \sim \pi_\phi(s)\}_{i=1}^m$ and $\{a_j \sim B(s)\}_{j=1}^n$, where $B(s)$ represents the set of actions appearing the the dataset B
- 8: Update ϕ, α by minimizing

$$\pi_\phi := \max_{\pi \in \Delta_{|S|}} \mathbb{E}_{s \sim B} \mathbb{E}_{a \sim \pi(\cdot|s)} \left[\min_{j=1, \dots, K} Q_{\theta_j}(s, a) \right] \quad \text{s.t.} \quad \mathbb{E}_{s \sim B} [\text{MMD}(B(s), \pi(\cdot|s))] \leq \epsilon, \quad (9)$$

where the MMD distance is defined as in (8)

- 9: Update Target Networks: $\theta'_i \leftarrow \tau \theta_i + (1 - \tau) \theta'_i, \phi' \leftarrow \tau \phi + (1 - \tau) \phi'$
 - 10: **end for**
- Ensure:** Policy π_ϕ and reward critics $\{Q_{\theta_j}\}$
-

507 **C.3 Detailed Analysis of Stitching Capability Verification**

508 We conduct a comprehensive evaluation of the cost and reward stitching abilities of the CQDT model.
 509 We summarize the parameters for testing the stitching ability in Table 6. To rigorously assess the cost
 510 stitching capabilities, we generate a suboptimal dataset by excluding the top $X\%$ of trajectories with
 511 the lowest cost returns from a set of trajectories that exhibit similar reward returns. Specifically, we
 512 group trajectories based on their RTG, then remove the trajectories with the lowest CTG in each group.
 513 In detail, we divide the trajectories into $\frac{\text{Max Return}}{10}$ groups, where Max Return denotes the highest
 514 return among all trajectories. Within each group, we remove trajectories with the lowest $X\%$ CTG.
 515 For the assessment of reward stitching capabilities, we similarly remove the top $X\%$ of trajectories
 516 with the highest reward returns from the trajectories. Figures 9 and 10 present visualizations of the
 517 datasets used to validate the reward and cost stitching abilities.

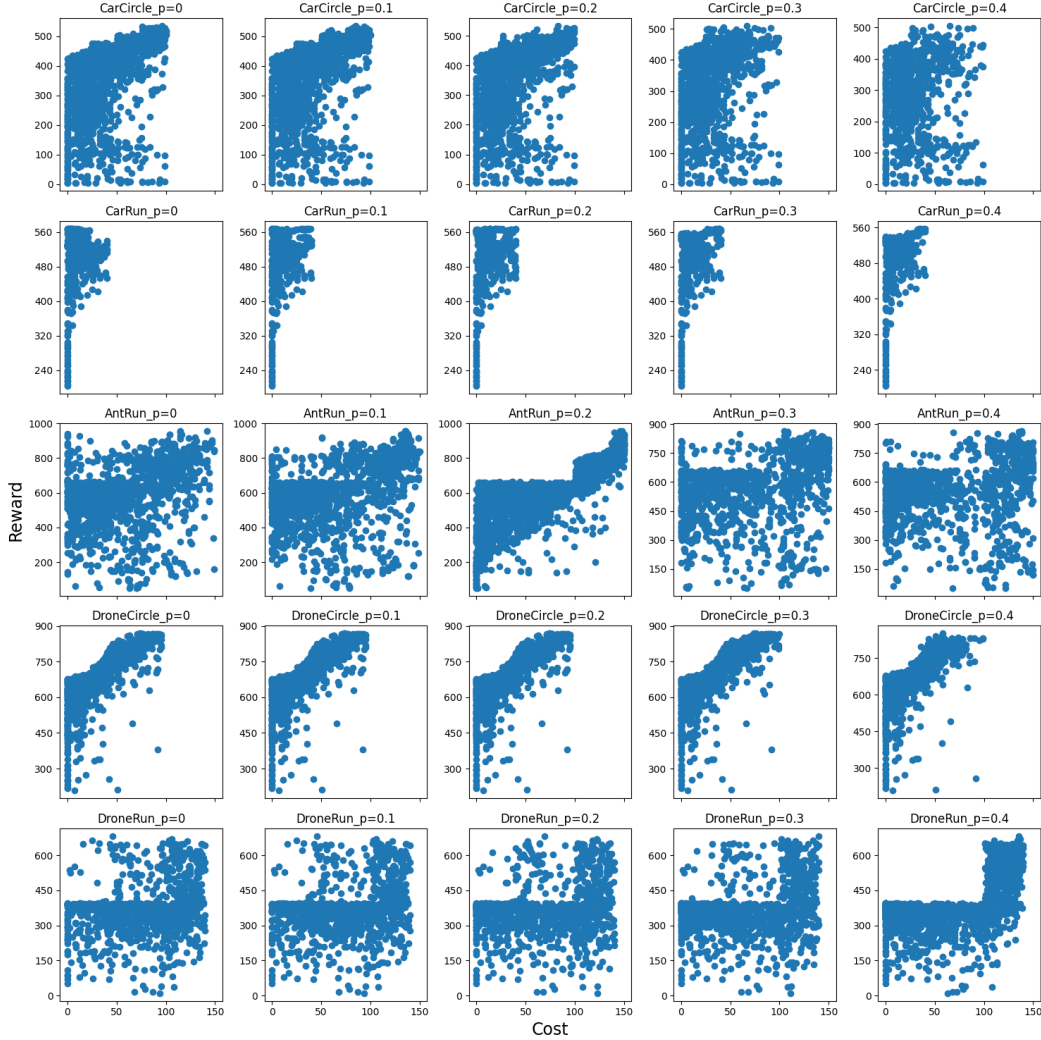


Figure 9: Visualization of the dataset used to validate reward stitching ability.

518 **C.4 Detailed Results from Zero-Shot Adaptation and Robustness Evaluation Experiment**

519 We study the robustness and the zero-shot adaptation ability of CQDT in this section. Figure 11
 520 depicts a fixed target cost, illustrating how variations in target return impact performance, providing
 521 a measure of robustness. Meanwhile, Figure 12 showcases a fixed target return, demonstrating
 522 how changes in target cost influence the actual performance, serving as an evaluation of zero-shot
 523 adaptation ability. For the raw experiment data, we include them in Table 8 and 9.

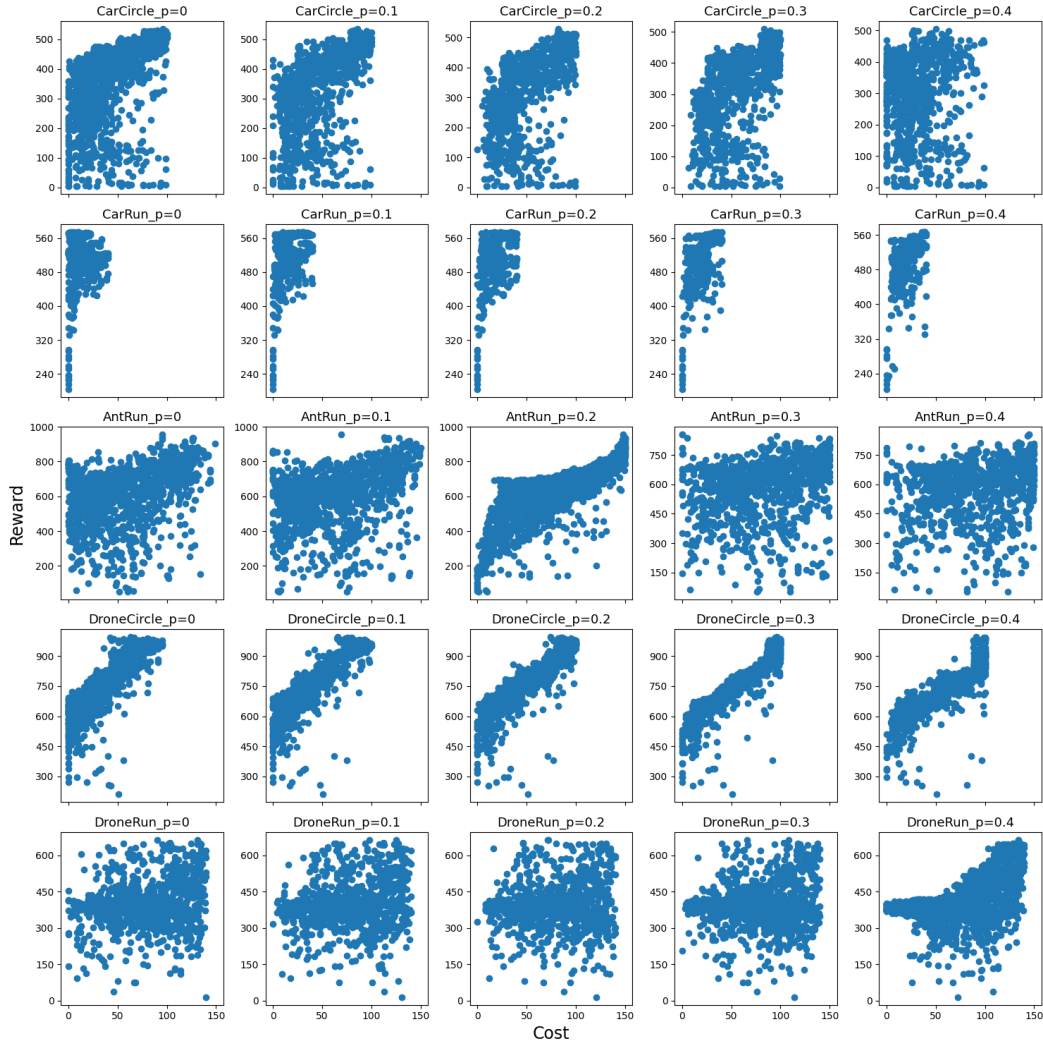


Figure 10: Visualization of the dataset used to validate cost stitching ability.

524 **Robustness Validation.** During the robustness validation phase, CQDT’s performance initially
 525 benefits from maintaining a constant target cost while gradually increasing the target return, leading
 526 to an augmentation in cumulative reward without violating cost constraints. However, once the target
 527 return reaches a certain threshold, further increases do not result in a proportional rise in cumulative
 528 reward. This phenomenon occurs because the target return setting guides CQDT’s prediction process
 529 but does not alter the intrinsic training process, which the model architecture and the training dataset
 530 determine. If the target return setting significantly exceeds the cumulative reward of the optimal
 531 trajectory in the training dataset, CQDT’s performance may not experience substantial improvements.
 532 Despite CQDT’s stitching ability, which allows cumulative rewards to increase when the target return
 533 setting exceeds the maximum threshold in the dataset, the positive effect of this ability is limited.

534 **Zero-shot Adaptation Validation.** In the zero-shot adaptation validation experiment, where the
 535 target return remains constant while the target cost varies, the analysis of cumulative rewards across
 536 five tasks reveals an initial increase followed by stabilization. As the target cost setting gradually
 537 increases, the constraints on cost become more relaxed, thereby enhancing the strategy’s ability to
 538 maximize cumulative rewards during the learning process. However, when the target cost setting
 539 exceeds the maximum cumulative cost threshold in the dataset, further increases in the target cost do
 540 not enhance the strategy’s ability to maximize cumulative rewards, leading to no further increase in
 541 cumulative rewards obtained during the test phase.



Figure 11: Results of robustness verification to different reward returns. Each column represents a task. The x-axis denotes the target return. The first row shows the evaluated reward, and the second row shows the evaluated cost, both under different target return. The dashed line represents the predefined cost limitation.

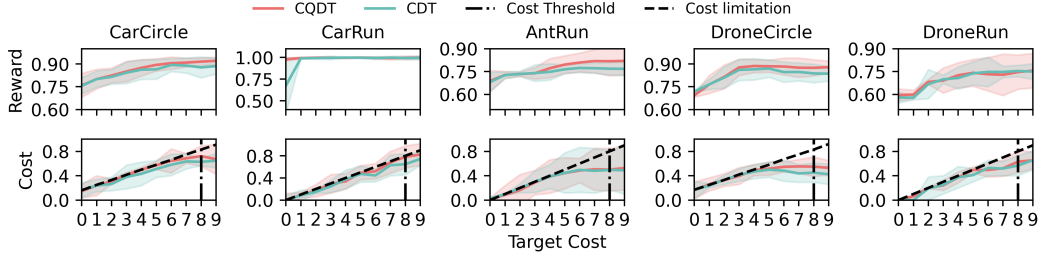


Figure 12: Results of zero-shot adaptation to different cost returns. Each column represents a task. The x-axis denotes the target cost return. The first row and the second row display the evaluated reward and cost under different target costs, respectively. The dashed line represents the predefined cost limitation, while the solid line indicates the maximum cost of trajectories included in the dataset.

Table 6: Parameter Settings for the Stitching Ability Validation Experiment

	CarCircle	CarRun	AntRun	DroneCircle	DroneRun
κ	90	35	130	90	120
Target Return	500	575	900	900	600
Target Cost	90	35	130	90	120

542 Indeed, while CQDT consistently outperforms CDT in both robustness and zero-shot adaptation
 543 experiments, the challenge of selecting appropriate target return and target cost values remains
 544 common to both approaches. Optimal choices for these parameters should align with the specific
 545 characteristics of the task training dataset. Setting a larger target return and a smaller target cost
 546 may not be advisable; a more viable approach involves tailoring these targets based on the inherent
 547 properties of the task-specific training data. This task-specific customization ensures a more effective
 548 and contextually appropriate utilization of the reinforcement learning framework.

Table 7: Settings of variables $cX-r[01-10]$ and $c[01-10]-rX$ in robustness validation and zero-shot adaptation experiments across different tasks. During the evaluation stage, the settings for Target Return and Target Cost correspond to the actual values in the environment, rather than the normalized values.

Task	Experiment	1	2	3	4	5	6	7	8	9	10
CarCircle	Robustness	c90r100	c90r150	c90r200	c90r250	c90r300	c90r350	c90r400	c90r450	c90r500	c90r550
	Zero-shot	c20r500	c30r500	c40r500	c50r500	c60r500	c70r500	c80r500	c90r500	c100r500	c110r500
CarRun	Robustness	c35r70	c35r140	c35r210	c35r280	c35r350	c35r420	c35r490	c35r560	c35r630	c35r700
	Zero-shot	c0r575	c5r575	c10r575	c15r575	c20r575	c25r575	c30r575	c35r575	c40r575	c45r575
AntRun	Robustness	c130r0	c130r100	c130r200	c130r300	c130r400	c130r500	c130r600	c130r700	c130r800	c130r900
	Zero-shot	c0r900	c20r900	c40r900	c60r900	c80r900	c100r900	c120r900	c140r900	c160r900	c180r900
DroneCircle	Robustness	c90r300	c90r400	c90r500	c90r600	c90r700	c90r800	c90r900	c90r1000	c90r1100	c90r1200
	Zero-shot	c20r900	c30r900	c40r900	c50r900	c60r900	c70r900	c80r900	c90r900	c100r900	c110r900
DroneRun	Robustness	c120r0	c120r100	c120r200	c120r300	c120r400	c120r500	c120r600	c120r700	c120r800	c120r900
	Zero-shot	c0r600	c20r600	c40r600	c60r600	c80r600	c100r600	c120r600	c140r600	c160r600	c180r600

Table 8: Experimental data from the CQDT, encompassing both the Zero-shot adaptation experiment (c[01-10]-rx) and the Robustness validation experiment (cx-r[01-10]) across diverse tasks. In the Robustness validation experiment, the target return is systematically varied while maintaining a constant target cost. Conversely, in the Zero-shot adaptation experiment, the target return is held constant while adjusting the size of the target cost. Selection criteria for cx-r[01-10] and c[01-10]-rx in different tasks are detailed for each scenario, with specific task names and corresponding the above table.

	CarCircle		CarRun		AntRun		DroneCircle		DroneRun	
	reward	cost	reward	cost	reward	cost	reward	cost	reward	cost
cX-r01	0.216 \pm 0.145	0.197 \pm 0.106	0.611 \pm 0.389	0.949 \pm 0.251	0.779 \pm 0.113	0.843 \pm 0.188	0.806 \pm 0.098	0.961 \pm 0.272	0.008 \pm 0.017	0.000 \pm 0.008
cX-r02	0.335 \pm 0.068	0.316 \pm 0.079	0.634 \pm 0.463	0.964 \pm 0.150	0.799 \pm 0.124	0.878 \pm 0.314	0.852 \pm 0.070	1.020 \pm 0.280	0.040 \pm 0.053	0.000 \pm 0.150
cX-r03	0.422 \pm 0.045	0.412 \pm 0.261	0.690 \pm 0.510	0.990 \pm 0.124	0.839 \pm 0.080	0.934 \pm 0.243	0.861 \pm 0.058	1.022 \pm 0.256	0.189 \pm 0.431	0.000 \pm 0.050
cX-r04	0.484 \pm 0.055	0.549 \pm 0.252	0.666 \pm 0.436	0.956 \pm 0.301	0.846 \pm 0.090	0.999 \pm 0.185	0.863 \pm 0.069	1.025 \pm 0.631	0.615 \pm 0.328	0.576 \pm 0.124
cX-r05	0.597 \pm 0.028	0.664 \pm 0.250	0.792 \pm 0.207	0.974 \pm 0.140	0.854 \pm 0.106	0.979 \pm 0.237	0.863 \pm 0.055	0.998 \pm 0.247	0.580 \pm 0.280	0.501 \pm 0.141
cX-r06	0.715 \pm 0.005	0.764 \pm 0.279	0.822 \pm 0.184	0.901 \pm 0.299	0.844 \pm 0.050	0.922 \pm 0.232	0.866 \pm 0.056	0.994 \pm 0.206	0.690 \pm 0.130	0.637 \pm 0.246
cX-r07	0.786 \pm 0.020	0.858 \pm 0.271	0.883 \pm 0.122	0.911 \pm 0.403	0.846 \pm 0.064	0.862 \pm 0.284	0.875 \pm 0.064	1.004 \pm 0.418	0.748 \pm 0.109	0.824 \pm 0.068
cX-r08	0.868 \pm 0.022	0.915 \pm 0.255	0.991 \pm 0.020	0.997 \pm 0.203	0.841 \pm 0.107	0.899 \pm 0.247	0.874 \pm 0.057	1.017 \pm 0.205	0.725 \pm 0.114	0.744 \pm 0.148
cX-r09	0.909 \pm 0.032	0.963 \pm 0.241	0.996 \pm 0.009	0.944 \pm 0.341	0.826 \pm 0.060	0.840 \pm 0.183	0.885 \pm 0.049	1.050 \pm 0.208	0.733 \pm 0.143	0.725 \pm 0.267
cX-r10	0.900 \pm 0.043	0.896 \pm 0.257	0.997 \pm 0.010	0.676 \pm 0.553	0.814 \pm 0.077	0.758 \pm 0.296	0.892 \pm 0.033	1.003 \pm 0.108	0.716 \pm 0.146	0.672 \pm 0.153
c01-rX	0.759 \pm 0.079	0.148 \pm 0.111	0.978 \pm 0.019	0.000 \pm 0.120	0.689 \pm 0.069	0.015 \pm 0.020	0.693 \pm 0.012	0.163 \pm 0.138	0.596 \pm 0.040	0.002 \pm 0.013
c02-rX	0.800 \pm 0.064	0.237 \pm 0.163	0.993 \pm 0.007	0.088 \pm 0.052	0.727 \pm 0.024	0.085 \pm 0.155	0.765 \pm 0.059	0.258 \pm 0.108	0.599 \pm 0.029	0.063 \pm 0.022
c03-rX	0.824 \pm 0.049	0.309 \pm 0.041	0.997 \pm 0.005	0.168 \pm 0.112	0.733 \pm 0.014	0.164 \pm 0.255	0.814 \pm 0.078	0.329 \pm 0.104	0.682 \pm 0.017	0.199 \pm 0.011
c04-rX	0.850 \pm 0.067	0.412 \pm 0.047	0.997 \pm 0.002	0.293 \pm 0.167	0.739 \pm 0.039	0.289 \pm 0.200	0.877 \pm 0.054	0.395 \pm 0.138	0.691 \pm 0.082	0.241 \pm 0.084
c05-rX	0.875 \pm 0.054	0.498 \pm 0.035	0.996 \pm 0.010	0.338 \pm 0.142	0.772 \pm 0.063	0.387 \pm 0.275	0.886 \pm 0.043	0.493 \pm 0.091	0.726 \pm 0.082	0.374 \pm 0.092
c06-rX	0.894 \pm 0.040	0.573 \pm 0.069	0.998 \pm 0.006	0.494 \pm 0.126	0.796 \pm 0.058	0.442 \pm 0.255	0.885 \pm 0.040	0.516 \pm 0.109	0.742 \pm 0.088	0.496 \pm 0.049
c07-rX	0.905 \pm 0.041	0.644 \pm 0.198	0.994 \pm 0.013	0.520 \pm 0.140	0.809 \pm 0.062	0.485 \pm 0.200	0.883 \pm 0.055	0.547 \pm 0.087	0.732 \pm 0.109	0.495 \pm 0.116
c08-rX	0.909 \pm 0.029	0.686 \pm 0.080	0.992 \pm 0.019	0.681 \pm 0.199	0.819 \pm 0.077	0.506 \pm 0.345	0.875 \pm 0.064	0.553 \pm 0.114	0.728 \pm 0.136	0.523 \pm 0.133
c09-rX	0.915 \pm 0.025	0.723 \pm 0.169	0.996 \pm 0.020	0.780 \pm 0.200	0.818 \pm 0.096	0.500 \pm 0.350	0.875 \pm 0.054	0.549 \pm 0.159	0.746 \pm 0.113	0.631 \pm 0.183
c10-rX	0.920 \pm 0.023	0.672 \pm 0.245	0.998 \pm 0.018	0.824 \pm 0.196	0.820 \pm 0.091	0.526 \pm 0.360	0.879 \pm 0.040	0.530 \pm 0.120	0.756 \pm 0.114	0.653 \pm 0.140

Table 9: Experimental data from the CDT, encompassing both the Zero-shot adaptation experiment (c[01-10]-rx) and the Robustness validation experiment (cx-r[01-10]) across diverse tasks. Other settings are the same as CQDT.

	CarCircle		CarRun		AntRun		DroneCircle		DroneRun	
	reward	cost	reward	cost	reward	cost	reward	cost	reward	cost
cX-r01	0.269 \pm 0.127	0.214 \pm 0.134	0.704 \pm 0.299	0.494 \pm 0.220	0.769 \pm 0.102	0.889 \pm 0.303	0.797 \pm 0.069	1.011 \pm 0.245	0.047 \pm 0.002	0.000 \pm 0.042
cX-r02	0.379 \pm 0.022	0.318 \pm 0.061	0.514 \pm 0.489	0.119 \pm 0.093	0.765 \pm 0.116	0.830 \pm 0.385	0.805 \pm 0.079	0.959 \pm 0.374	0.071 \pm 0.016	0.000 \pm 0.158
cX-r03	0.453 \pm 0.026	0.352 \pm 0.086	0.788 \pm 0.214	0.741 \pm 0.116	0.803 \pm 0.073	0.913 \pm 0.349	0.831 \pm 0.075	1.016 \pm 0.484	0.162 \pm 0.189	0.000 \pm 0.017
cX-r04	0.528 \pm 0.009	0.507 \pm 0.166	0.990 \pm 0.014	0.983 \pm 0.274	0.798 \pm 0.089	0.872 \pm 0.359	0.841 \pm 0.077	1.009 \pm 0.258	0.546 \pm 0.195	0.671 \pm 0.137
cX-r05	0.609 \pm 0.031	0.575 \pm 0.193	0.991 \pm 0.002	0.999 \pm 0.259	0.782 \pm 0.089	0.860 \pm 0.310	0.840 \pm 0.049	0.985 \pm 0.126	0.640 \pm 0.343	0.714 \pm 0.297
cX-r06	0.700 \pm 0.014	0.682 \pm 0.188	0.988 \pm 0.014	0.979 \pm 0.250	0.829 \pm 0.055	0.976 \pm 0.224	0.839 \pm 0.064	1.021 \pm 0.168	0.723 \pm 0.204	0.701 \pm 0.171
cX-r07	0.785 \pm 0.029	0.773 \pm 0.146	0.988 \pm 0.021	0.949 \pm 0.309	0.813 \pm 0.054	0.868 \pm 0.194	0.847 \pm 0.080	0.983 \pm 0.328	0.732 \pm 0.097	0.721 \pm 0.245
cX-r08	0.863 \pm 0.031	0.847 \pm 0.184	0.992 \pm 0.017	0.966 \pm 0.206	0.776 \pm 0.076	0.698 \pm 0.278	0.851 \pm 0.060	1.009 \pm 0.113	0.717 \pm 0.068	0.764 \pm 0.303
cX-r09	0.889 \pm 0.041	0.838 \pm 0.164	0.996 \pm 0.009	0.957 \pm 0.443	0.805 \pm 0.030	0.750 \pm 0.135	0.856 \pm 0.057	1.022 \pm 0.084	0.729 \pm 0.063	0.694 \pm 0.298
cX-r10	0.874 \pm 0.054	0.862 \pm 0.251	0.993 \pm 0.014	0.956 \pm 0.587	0.777 \pm 0.028	0.638 \pm 0.147	0.864 \pm 0.043	1.023 \pm 0.213	0.731 \pm 0.082	0.746 \pm 0.300
c01-rX	0.756 \pm 0.054	0.161 \pm 0.123	0.675 \pm 0.322	0.000 \pm 0.080	0.681 \pm 0.064	0.029 \pm 0.010	0.718 \pm 0.007	0.165 \pm 0.144	0.579 \pm 0.036	0.003 \pm 0.003
c02-rX	0.800 \pm 0.054	0.239 \pm 0.103	0.990 \pm 0.006	0.083 \pm 0.097	0.726 \pm 0.022	0.095 \pm 0.014	0.765 \pm 0.033	0.238 \pm 0.037	0.579 \pm 0.014	0.000 \pm 0.000
c03-rX	0.815 \pm 0.049	0.264 \pm 0.111	0.992 \pm 0.011	0.148 \pm 0.032	0.734 \pm 0.017	0.194 \pm 0.073	0.810 \pm 0.070	0.332 \pm 0.076	0.668 \pm 0.095	0.191 \pm 0.191
c04-rX	0.843 \pm 0.071	0.380 \pm 0.195	0.995 \pm 0.020	0.247 \pm 0.113	0.740 \pm 0.054	0.301 \pm 0.178	0.861 \pm 0.069	0.413 \pm 0.104	0.699 \pm 0.028	0.250 \pm 0.150
c05-rX	0.862 \pm 0.064	0.431 \pm 0.194	0.998 \pm 0.006	0.373 \pm 0.167	0.748 \pm 0.054	0.387 \pm 0.198	0.864 \pm 0.070	0.468 \pm 0.115	0.707 \pm 0.073	0.375 \pm 0.175
c06-rX	0.864 \pm 0.065	0.511 \pm 0.139	0.997 \pm 0.006	0.472 \pm 0.148	0.765 \pm 0.029	0.446 \pm 0.330	0.870 \pm 0.058	0.499 \pm 0.126	0.738 \pm 0.077	0.411 \pm 0.111
c07-rX	0.893 \pm 0.043	0.580 \pm 0.178	0.993 \pm 0.013	0.445 \pm 0.115	0.773 \pm 0.029	0.497 \pm 0.370	0.846 \pm 0.066	0.483 \pm 0.175	0.748 \pm 0.097	0.528 \pm 0.128
c08-rX	0.889 \pm 0.049	0.635 \pm 0.123	0.997 \pm 0.011	0.630 \pm 0.190	0.772 \pm 0.034	0.481 \pm 0.367	0.847 \pm 0.059	0.438 \pm 0.146	0.747 \pm 0.038	0.518 \pm 0.118
c09-rX	0.877 \pm 0.069	0.628 \pm 0.063	0.993 \pm 0.017	0.648 \pm 0.192	0.769 \pm 0.041	0.495 \pm 0.326	0.837 \pm 0.061	0.451 \pm 0.158	0.754 \pm 0.037	0.566 \pm 0.167
c10-rX	0.886 \pm 0.052	0.646 \pm 0.120	0.995 \pm 0.021	0.743 \pm 0.117	0.768 \pm 0.045	0.490 \pm 0.392	0.836 \pm 0.049	0.424 \pm 0.167	0.750 \pm 0.051	0.655 \pm 0.155