# PRETRAINING A SHARED Q-NETWORK FOR DATA EFFICIENT OFFLINE REINFORCEMENT LEARNING

Anonymous authors

Paper under double-blind review

# ABSTRACT

Offline reinforcement learning (RL) aims to learn a policy from a static dataset without further interactions with the environment. Collecting sufficiently large datasets for offline RL is exhausting since this data collection requires colossus interactions with environments and becomes tricky when the interaction with the environment is restricted. Hence, how an agent learns the best policy with a minimal static dataset is a crucial issue in offline RL, similar to the sample efficiency problem in online RL. In this paper, we propose a simple yet effective plug-andplay pretraining method to initialize a feature of a Q-network to enhance data efficiency in offline RL. Specifically, we introduce a shared Q-network structure that outputs predictions of the next state and Q-value. We pretrain the shared Q-network through a supervised regression task that predicts a next state and trains the shared Q-network using diverse offline RL methods. Through extensive experiments, we empirically demonstrate that the proposed method enhances the performance of existing popular offline RL methods on the D4RL and Robomimic benchmarks, with an average improvement of 135.94% on the D4RL benchmark. Furthermore, we show that the proposed method significantly boosts data-efficient offline RL across various data qualities and data distributions. Notably, our method adapted with only 10% of the dataset outperforms standard algorithms even with full datasets.

028 029

031

004

006

008 009

010 011

012

013

014

015

016

017

018

019

021

024

025

026

027

## 1 INTRODUCTION

Sample efficiency is a crucial issue in reinforcement learning (RL) since typical RL considers an online learning nature that involves iterative processes between experience collections and policy improvements through online interactions with the environment (Sutton et al., 1998). Unfortunately, requiring excessive online interactions is impractical in several cases since data collection requires expensive costs and retains potential risks of the agent, e.g. hardware corruption. Offline RL is one approach to alleviate this sample efficiency problem, which provides a solution by avoiding online interactions with the environment (Levine et al., 2020). In recent years, pretraining with offline RL and fine-tuning with online RL have been investigated to improve sample efficiency of the online interactions (Nakamoto et al., 2024; Xie et al., 2021; Rafailov et al., 2023; Ball et al., 2023).

041 Similar to addressing the sample efficiency problem in online RL, learning offline RL with minimal 042 datasets is necessary since collecting enormous experience charges expensive costs and unfavorable 043 explorations, hampering the possibility of offline RL in the real world. In this paper, we name this 044 problem as data efficiency where an agent tries to learn the best policy with minimal data in the offline RL scheme. Despite the necessity of data efficiency, this problem has not been treated enough in previous works. Although some researchers have evaluated their work empirically on reduced datasets 046 in part of the experiments (Agarwal et al., 2020; Kumar et al., 2020a;b), they have overlooked this 047 data efficiency problem. In the case of online RL, model-based RL and representation learning have 048 proposed the resolution of sample efficiency problem (Sutton, 1991; Hafner et al., 2019b; Schwarzer et al., 2020; 2021). As in online RL, one can expect that offline model-based RL or representation method might resolve this data efficiency problem (Yu et al., 2020; Sun et al., 2023; Yang & Nachum, 051 2021). However, Figure 7 demonstrates that both approaches are unable to overcome this problem. 052

In this work, we propose a simple yet effective plug-and-play method that pretrains a shared Qnetwork toward data-efficient offline RL. Specifically, the shared Q-network structure is composed of



Figure 1: **Overview of our pretraining method.** Our method splits the original Q-network into two core architectures: a shared network that extracts the representation z from the concatenated vector of state s and action a and separated heads for training the transition model network and Q-network, respectively.

two parts as illustrated in Figure 1. First, a shared deep neural network layer  $(h_{\varphi})$  takes the state and action pair as inputs. Second, separate shallow output parts  $(g_{\psi} \text{ and } f_{\theta})$  consist of two linear layers that individually output a Q-value for a Q-function and a next state prediction for a transition model. The learning phase of the shared Q-network consists of a pretraining and an RL training phase. In the pretraining phase, the shared network attached with a shallow transition layer  $(h_{\varphi} \text{ and } g_{\psi})$  is trained through a supervised regression task that predicts the transition model. After the pretraining phase where the shared network is initialized with the pretraining, the shared network is connected with a shallow Q layer  $(h_{\varphi} \text{ and } f_{\theta})$  and trained with an existing offline RL value learning.

We empirically demonstrate that our method improves the performance of existing popular offline RL 077 methods on the D4RL (Fu et al., 2020), and Robomimic (Mandlekar et al., 2021), benchmarks with an average improvement of 135.94% on the D4RL benchmark. We also show that our method maintains 079 data-efficient performance with fragments of the dataset across the *data quality* on the D4RL dataset. Moreover, we investigate our method across the *data collection strategies* on the ExoRL datasets 081 (Yarats et al., 2022), assuming a small dataset would have a shifted data distribution compared to a 082 large dataset. As a result, we demonstrate that our method improves the performance regardless of 083 the qualities of the datasets and the data distributions. Figure 6 and Figure 9 show that our method 084 with 10% of datasets outperforms vanilla algorithms even with full datasets. Furthermore, Figure 7 085 demonstrates that our method indeed outperforms the offline model-based RL and representation approaches in reduced datasets.

087

090

# 2 RELATED WORKS

Offline RL. Offline RL aims to learn a policy with static data without further interactions with the 091 environment. Previous approaches have mainly addressed the distribution shift problem, which is 092 caused by the idea that queries of the Q-function over out-of-distribution actions may yield overly optimistic values during offline training (Fujimoto et al., 2019; Kumar et al., 2019; Levine et al., 094 2020; Kumar et al., 2020b; Fujimoto & Gu, 2021; Kostrikov et al., 2021a). Recently, scalability to a large dataset and neural network model has been studied (Chebotar et al., 2023; Padalkar et al., 096 2023; Team et al., 2024). In other fields, pretraining with offline RL and fine-tuning with online RL is examined to improve sample efficiency in the online interaction step (Nakamoto et al., 2024; Xie 098 et al., 2021; Rafailov et al., 2023; ?). In contrast, distinct experiments over the way to consuming 099 the static dataset have been conducted, e.g., an imbalanced dataset, unlabeled data, and even data corruption under an offline RL scheme (Hong et al., 2023; Yu et al., 2022; Yang et al., 2023). While 100 prior research (Agarwal et al., 2020; Kumar et al., 2020a;b) often has evaluated their work on reduced 101 datasets as a partial result, the field overlooks the data efficiency problem itself as a main contribution. 102 In contrast, we aim to improve the data efficiency in offline RL (i.e., learning the best policy with 103 minimal data). In this work, we propose a simple yet effective plug-and-play method for pretraining 104 a shared Q-network toward the data-efficient offline RL. 105

Sample efficient RL. A common issue in most RL algorithms is sample efficiency: excessive interactions with the environment are required to learn an optimal policy. For this reason, sample efficiency has been an active research topic in RL (Kostrikov et al., 2021b; Yarats et al., 2021c;

108 D'Oro et al., 2022). Model-based RL (Sutton, 1991; Deisenroth & Rasmussen, 2011; Hafner et al., 109 2019b;a; Hansen et al., 2022) is a common approach to resolve sample inefficiency by learning a 110 (latent) dynamics model and using it to generate additional transition samples. Otherwise, effective 111 pretraining (Schwarzer et al., 2021; Yarats et al., 2021c) and data augmentation (Laskin et al., 2020; 112 Kostrikov et al., 2021b) play a critical role in improving sample efficiency in RL. Recently, offline-toonline (Lee et al., 2022; Ball et al., 2023; Rafailov et al., 2023; Feng et al., 2023; Nakamoto et al., 113 2024) and foundation model (Ahn et al., 2022; Seo et al., 2022; Brohan et al., 2023b;a; Bhateja et al., 114 2023) have tackled this problem where the poor sample efficiency of online RL regime is alleviated 115 by leveraging large offline data. In this paper, we separately define the data efficiency problem in 116 offline RL as the ability of an offline RL algorithm how an agent can learn the best policy even with 117 minimal pre-collected samples called dataset in offline RL. We claim that this data efficiency problem 118 is different from the sample efficiency problem since online RL has opportunities for interactions 119 with environments which can present another chance to improve the sample efficiency. 120

121

# **3** MARKOV DECISION PROCESS

122 123

135

136

144 145

146

151

154

We consider the Markov decision process, where the agent sequentially takes actions to maximize cumulative discounted rewards. In a Markov decision process with the state-space  $S := \{1, 2, ..., |S|\}$ and action-space  $A := \{1, 2, ..., |A|\}$ , the decision maker selects an action  $a \in A$  at the current state  $s \in S$ , then the state transits to the next state  $s' \in S$  with probability P(s'|s, a), and the transition incurs a reward  $r(s, a, s') \in \mathbb{R}$ , where P(s'|s, a) is the state transition probability from the current state  $s \in S$  to the next state  $s' \in S$  under action  $a \in A$ , and r(s, a, s') is the reward function. For convenience, we consider a deterministic reward function and simply write  $r(s_k, a_k, s_{k+1}) =: r_k, k \in \{0, 1, ...\}$ .

A deterministic policy,  $\pi : S \to A$ , maps a state  $s \in S$  to an action  $\pi(s) \in A$ . The objective of the Markov decision problem is to find a deterministic (or stochastic) optimal policy,  $\pi^*$ , such that the cumulative discounted rewards over infinite time horizons is maximized, i.e.,

$$\pi^* := \operatorname*{arg\,max}_{\pi \in \Theta} \mathbb{E}\left[ \left| \sum_{k=0}^{\infty} \gamma^k r_k \right| \pi \right].$$

where  $\gamma \in [0, 1)$  is the discount factor,  $\Theta$  is the set of all deterministic policies,  $(s_0, a_0, s_1, a_1, ...)$  is a state-action trajectory generated by the Markov chain under policy  $\pi$ , and  $\mathbb{E}[\cdot|\pi]$  is an expectation conditioned on the policy  $\pi$ . Moreover, Q-function under policy  $\pi$  is defined as

$$Q^{\pi}(s,a) = \mathbb{E}\left[\left|\sum_{k=0}^{\infty} \gamma^{k} r_{k}\right| s_{0} = s, a_{0} = a, \pi\right], \quad (s,a) \in \mathcal{S} \times \mathcal{A}$$

# 4 PRETRAINING Q-NETWORK WITH TRANSITION MODEL HELPS IMPROVING DATA EFFICIENCY

In this paper, we propose a simple yet effective pretraining method adapting features of the transition
model into the initialization of Q-network to improve data efficiency in offline RL. To this end, we
design Q-network that partially shares a network with the estimation of the transition model. In
particular, the transition model is constructed as follows:

$$\hat{s}' = (g_{\psi} \circ h_{\varphi})(s, a), \quad (s, a) \in \mathcal{S} \times \mathcal{A}, \tag{1}$$

where  $\hat{s}'$  is the estimated next state,  $g_{\psi}$  is a parameterized linear function, and  $h_{\varphi}$  is shared with the *Q*-network, which is defined as

$$Q_{\varphi,\theta}(s,a) = (f_{\theta} \circ h_{\varphi})(s,a), \quad (s,a) \in \mathcal{S} \times \mathcal{A},$$
(2)

where  $f_{\theta}$  is also a parameterized linear function that represents the linear output layer and  $h_{\varphi}$ represents the fully connected neural network layers shared with the transition model in (1). The overall structures of the neural networks are illustrated in Figure 1.

In the proposed method, the transition model  $g_{\psi} \circ h_{\varphi}$  is pretrained by minimizing the mean squared prediction error loss function

160  
161
$$\mathcal{L}_{pre}(\varphi,\psi) = \sum_{(s,a,s')\in\mathcal{D}} (s' - (g_{\psi} \circ h_{\varphi})(s,a))^2$$
(3)

162 Algorithm 1 Pretraining Q-network scheme for Offline RL 163 **Input**: Dataset  $\mathcal{D}$  of transition (s, a, s'), learning rate  $\alpha$ 164 Initialize parameters  $\varphi, \psi$ for each gradient step do 166 Sample a mini-batch  $\mathcal{B} \sim \mathcal{D}$ 167 Compute the transition model estimation error  $\mathcal{L}_{pre}(\varphi,\psi) = \sum_{(s,a,s')\in\mathcal{B}} (s' - (g_{\psi} \circ h_{\varphi})(s,a))^2$ 169 170 171 Update weights of the shared network and transition model network 172  $\varphi \leftarrow \varphi - \alpha \nabla_{\varphi} \mathcal{L}_{pre}(\varphi, \psi), \quad \psi \leftarrow \psi - \alpha \nabla_{\psi} \mathcal{L}_{pre}(\varphi, \psi)$ 173 174 end for 175 **Output**: Pretrained weights  $\varphi$  of the shared network 176 177 178 over the pre-collected dataset  $\mathcal{D}$  which includes a given set of the transition (s, a, s'). Afterward, 179 the pretrained parameter  $\varphi$  can be used as an initial or fixed parameter for standard RL algorithms based on the Q-network structure in (4) without any modification. The overall pretraining process is 181 summarized in Algorithm 1 for offline RL. We also note that similar principles can be applied for

Later in this paper, we empirically demonstrate that combining the proposed pretraining method with
 existing offline RL methods can effectively improve their performances. Moreover, we demonstrate
 that our method indeed improves data efficiency through some experiment settings in offline RL.

### 4.1 ANALYSIS: BASED ON THE PROJECTED BELLMAN EQUATION

182

186 187

188

195

196

202

203 204 205

210

040

online RL as well, and the corresponding algorithm is given in Appendix A.

In this section, we analyze how our method can resolve the data efficiency problem from the perspective of the projected Bellman equation. For simplicity and convenience of presentation, we assume that the state and action spaces are discrete and finite, and the transition is deterministic. However, the principles in this paper can be extended to more general continuous state and continuous action cases. Our analysis is based on the observation that Q-function with neural networks can be generally represented by (2). Defining the feature vector  $z = h_{\varphi}(s, a) \in \mathbb{R}^m$ , it can be rewritten as

$$Q_{\varphi,\theta}(s,a) = \sum_{i=1}^{m} \theta_i h_{\varphi,i}(s,a) = \langle \theta, h_{\varphi}(s,a) \rangle, \quad (s,a) \in \mathcal{S} \times \mathcal{A}.$$
(4)

197 198 When  $\varphi$  is fixed, then the above structure can be viewed as a linear function approximation with the feature function  $h_{\varphi,i}$ . In the proposed method,  $h_{\varphi,i}$  is indeed pretrained by minimizing the loss in (3) and then fixed while learning Q-function in (4). Therefore, the interpretation based on the linear function approximation is expected to be a reasonable model to explain the phenomenon in the proposed method.

It is well known that with linear function approximation, the corresponding standard Bellman equation

$$Q_{\varphi,\theta}(s,a) = R(s,a) + \gamma \sum_{s' \in S} P^{\pi}(s'|s,a) \sum_{a' \in A} Q_{\varphi,\theta}(s',a')$$

may not admit a solution in general. However, typical TD-learning algorithms are known to converge to the unique fixed point of the projected Bellman equation. In particular, considering the vector form of the Bellman equation,  $Q_{\varphi,\theta} = R + \gamma P^{\pi} Q_{\varphi,\theta}$ , the projected Bellman equation (Melo & Ribeiro (2007)) is known to admit a solution

$$Q_{\varphi,\theta} = \Pi(R + \gamma P^{\pi} Q_{\varphi,\theta})$$

where  $\Pi$  is the projection onto the column space,  $C(H_{\varphi})$ , of the feature matrix  $H_{\varphi}$  defined as

212  
213  
214  

$$H_{\varphi} := \begin{bmatrix} \vdots \\ h_{\varphi}(s,a)^T \end{bmatrix}.$$

Table 1: The Rank of the latent space of Q-network on the D4RL benchmark. We compare the
 rank of the latent space between a vanilla TD3+BC and TD3+BC adapted with our method over 512
 samples. As a result, adapting our method significantly increases the rank of the latent space, leading
 to reduced approximation error.

	Halfcheatah			Hopper	Walker2d		
	TD3+BC	TD3+BC (+ours)	TD3+BC	TD3+BC (+ours)	TD3+BC	TD3+BC (+ours)	
Random	59	236	69	192	72	82	
Medium	55	249	85	227	55	254	
Medium Replay	49	252	77	249	77	255	
Medium Expert	58	236	86	232	52	253	
Expert	44	198	104	198	68	225	



Figure 2: Reduced approximation error with the expanded column space of  $H_{\varphi}$ . In linear approximation, there exists  $Q^{\pi}$  outside of the column space of  $H_{\varphi}$ . To deal with this problem, the projected Bellman equation projects  $Q^{\pi}$  to  $\Pi Q^{\pi}$  which exists in the column space of  $H_{\varphi}$ .

243 The corresponding solution is known to have the error bound

244 245

235 236 237

238

239

240 241 242

220 221 222

246

247 248  $||Q_{\varphi,\theta} - Q^{\pi}||_{\infty} \le \frac{1}{1-\gamma} ||\Pi Q^{\pi} - Q^{\pi}||_{\infty},$  (5)

where  $Q^{\pi}$  is the true Q-function corresponding to the target policy  $\pi$ . As can be seen from the above bound, the error depends on the feature matrix  $H_{\varphi}$ . We can observe that the smaller the distance between  $C(H_{\varphi})$  and  $Q^{\pi}$ , the smaller the error between  $Q_{\varphi,\theta}$  and  $Q^{\pi}$ . Therefore, a proper choice of the feature function is key to the successful estimation of  $Q^{\pi}$ .

With the neural network function approximation, typical value-based RL algorithms update both  $\varphi$ and  $\theta$  simultaneously via TD-learning algorithms. Since the feature functions,  $h_{\varphi,i}$ , are in general nonlinear and non-convex in  $\varphi$ , it may sometimes converge to a local optimal solution. This in turn implies that appropriate initialization or pretraining of the feature functions,  $h_{\varphi,i}$ , can play an important role for estimating *Q*-function with smaller approximation errors on the right-hand side of (5) by avoiding suboptimal local solutions.

We conjecture that the pretraining approach with the transition model introduced in the previous section can effectively shape the feature functions so that the column space  $C(H_{\varphi})$  can cover higher dimensional vector space in  $\mathbb{R}^{|S \times A|}$ . As shown in Figure 2, this eventually results in a reduction of the solution error on the right-hand side of (5). To support this, we empirically compare the rank of the Q-network in the latent space between vanilla and the pretrained TD3+BC with our method over 512 data samples.

Table 1 exhibits that adapting our method shows a significantly higher rank than the rank of the vanilla method. From the results, we claim that the proposed method indeed expands the column space  $C(H_{\varphi})$  and covers higher dimensional vector space in  $\mathbb{R}^{|S \times A|}$ , leading to more precise Q-function estimation. In other words, we might learn a more precise Q-function with the same amount of samples, and it means that we can get a desirably estimated Q-function with less data. In the following section, we demonstrate our claim with empirical experiments.

#### 270 5 **EXPERIMENTS** 271

272

273 In this section, we evaluate our method over existing offline RL methods with the popular offline RL benchmarks, D4RL, and the more complex domain, Robomimic. Furthermore, we examine the 274 proposed method over the partial fragments of D4RL and ExoRL datasets for data-efficient offline RL. 275 We introduce a detailed experimental setup and baselines in the following paragraphs and provide 276 empirical results subsequently. 277

278 **Experimental setup.** We have considered heterogeneous tasks and diverse datasets for precise 279 comparisons. For the locomotion task, the proposed method is compared with existing methods in the 280 popular D4RL benchmark (Fu et al., 2020). Three different embodied agents and five distinct datasets are considered in order to validate the effectiveness of the proposed method: HalfCheetah, Hopper, 281 Walker2d for agents and random, medium-replay, medium, medium-expert, expert for datasets. For the 282 tabletop manipulation tasks, we have evaluated the proposed method in the Robomimic benchmark, 283 (Mandlekar et al., 2021), where off-the-shelf offline RL methods are already implemented. Two 284 different tabletop tasks and mixed-quality datasets are considered to verify the scalability of the 285 proposed method: Lift, Can for tasks and Machine-Generated (MG) for datasets. For data-efficient 286 offline RL, we have evaluated the proposed method across the reward qualities of the datasets of 287 D4RL Gym locomotion tasks, and the dataset collection strategies for walker walk (i.e. SMM, RND, 288 ICM) and point mass maze (i.e. Proto, Diayn) in ExoRL (Yarats et al., 2022). See Appendix C for a 289 more detailed setup for tasks and datasets.

290 Baselines. We have designed extensive experiments on the D4RL benchmark to verify the effective-291 ness of the proposed method built on top of the popular offline RL methods, including AWAC (Nair 292 et al., 2020), CQL (Kumar et al., 2020b), TD3+BC (Fujimoto & Gu, 2021), and IQL (Kostrikov 293 et al., 2021a). To verify the benefits of the proposed method, we compared the normalized scores between the vanilla method and the one combined with the proposed pretraining method. Similar 295 to the D4RL benchmark, the success rate is compared on the Robomimic benchmark, where IQL, 296 TD3+BC, BCQ (Fujimoto et al., 2019), and IRIS (Mandlekar et al., 2020), were used in combination 297 with the proposed methods. We also evaluate MOPO (Yu et al., 2020), MOBILE (Sun et al., 2023) and ACL (Yang & Nachum, 2021) to compare the proposed method with offline model-based RL 298 and representation approaches. On the ExoRL benchmark, we used TD3 (Fujimoto et al., 2018), for 299 walker walk task, and CQL for point mass maze tasks. See Appendix E for more implementation 300 details. 301

302 303

304

305

306

307

308

309

Table 2: Averaged normalized scores on the D4RL benchmark over 5 seeds. In each column corresponding to different RL methods, values on the left-hand side are scores of the baseline methods directly taken from the literature. The values on the right-hand side of each column represent scores of the proposed methods combined with the baselines. The increased scores compared to the baselines are highlighted in blue font, and they are reported with the mean and standard deviations over five random seeds.

		AWAC	CQL	IQL	TD3+BC
Dondom	HalfCheetah	$2.2 \rightarrow 51.10 \pm 0.89$	$21.7 \pm 0.9 \rightarrow 31.94 \pm 2.63$	$\rightarrow$ 18.28 $\pm$ 1.02	$10.2 \pm 1.3 \rightarrow 14.83 \pm 0.54$
Kandoin	Walker2d	$5.1 \rightarrow 13.11 \pm 3.91$	$10.7 \pm 0.1 \rightarrow 30.20 \pm 2.00$ $2.7 \pm 1.2 \rightarrow 19.56 \pm 4.49$	$\rightarrow$ 10.07 $\pm$ 0.41 $\rightarrow$ 8.88 $\pm$ 0.71	$1.4 \pm 1.6 \rightarrow 11.23 \pm 5.05$
Madium	HalfCheetah	$37.4 \rightarrow 54.63 \pm 1.45$	$37.2 \pm 0.3 \rightarrow 39.93 \pm 18.84$	$47.4 \rightarrow 48.85 \pm 0.16$	$42.8 \pm 0.3 \rightarrow 49.17 \pm 0.26$
Wedium	Walker2d	$72.0 \rightarrow 101.75 \pm 0.20$ $30.1 \rightarrow 89.51 \pm 0.88$	$44.2\pm10.8 \rightarrow 90.38\pm2.23$ $57.5\pm8.3 \rightarrow 84.66\pm0.67$	$78.3 \rightarrow 83.63 \pm 1.14$	$99.5 \pm 1.0 \rightarrow 71.52 \pm 2.16$ $79.7 \pm 1.8 \rightarrow 87.09 \pm 0.60$
	HalfCheetah	$\rightarrow$ 55.75 $\pm$ 1.30	41.9±1.1→47.60±0.37	44.2→45.48±0.17	$43.3 \pm 0.5 \rightarrow 45.84 \pm 0.26$
Medium Replay	Hopper Walker2d	$\rightarrow 106.67 \pm 0.59$ $\rightarrow 100.31 \pm 2.11$	$28.6 \pm 0.9 \rightarrow 98.63 \pm 2.12$ $15.8 \pm 2.6 \rightarrow 87.66 \pm 1.30$	$94.7 \rightarrow 99.43 \pm 1.71$ $73.9 \rightarrow 87.95 \pm 1.68$	$31.4 \pm 3.0 \rightarrow 100.16 \pm 1.60$ $25.2 \pm 5.1 \rightarrow 92.01 \pm 1.58$
Malian Francis	HalfCheetah	$36.8 \rightarrow 90.05 \pm 1.89$	$27.1 \pm 3.9 \rightarrow 82.75 \pm 6.51$	$86.7 \rightarrow 95.25 \pm 0.14$	$97.9 \pm 4.4 \rightarrow 96.89 \pm 0.92$
Medium Expert	Walker2d	$80.9 \rightarrow 113.23 \pm 0.22$ $42.7 \rightarrow 111.88 \pm 0.28$	$68.1 \pm 13.1 \rightarrow 91.63 \pm 42.48$	$91.5 \rightarrow 105.77 \pm 11.31$ $109.6 \rightarrow 112.09 \pm 0.93$	$112.2 \pm 0.2 \rightarrow 113.02 \pm 0.19$ $101.1 \pm 9.3 \rightarrow 111.58 \pm 0.35$
	HalfCheetah	78.5→93.48±0.11	82.4±7.4→97.09±1.03	$\rightarrow$ 97.40 $\pm$ 0.13	105.7±1.9→98.86±0.55
Expert	Hopper Walker2d	$57.0 \rightarrow 111.22 \pm 0.35$	$111.2\pm2.1 \rightarrow 112.10\pm0.35$ $103.8\pm7.6 \rightarrow 110.64\pm0.28$	$\rightarrow$ 113.34 $\pm$ 0.46 $\rightarrow$ 112.80 $\pm$ 1.08	$112.2 \pm 0.2 \rightarrow 113.35 \pm 0.28$ $105.7 \pm 2.7 \rightarrow 111.00 \pm 0.15$
Total		$\rightarrow$ 1265.01 $\pm$ 48.07	764.3±61.5→1136.03±86.78	$\rightarrow$ 1118.46 $\pm$ 23.25	979.3±33.4→1148.12±14.65



Figure 3: Learning curves of TD3+BC. The blue and orange curves are, respectively, the normalized scores of TD3+BC and TD3+BC pretrained with the proposed method. The vertical red reference lines split the pretraining and main training phases. After the pretraining phase, TD3+BC combined with the proposed method quickly outperforms the vanilla TD3+BC by a large margin.



Figure 4: Averaged normalized scores across pretraining time-step rates. R, M, MR, ME, and E represent random, medium, medium replay, medium expert, and expert datasets on the D4RL benchmark, respectively.

#### 5.1PERFORMANCE IMPROVEMENT IN OFFLINE RL BENCHMARKS

To demonstrate the effectiveness of the proposed method over existing offline RL methods, we 356 evaluate our method on D4RL and Robomimic datasets. In Table 2, the normalized scores between 357 the vanilla and the one combined with our method are compared for each environment and dataset 358 in D4RL. One can observe that the proposed method combined with the baselines improves the corresponding original methods, achieving an average improvement of 135.94%, across diverse envi-359 ronments and datasets. Specifically, one can observe that all methods including AWAC (+306.45%), 360 CQL (+132.77%), IQL (+9.21%), and TD3+BC (+95.34%) exhibit significantly increased per-361 formance on average compared to the results reported in the original papers. We have taken all 362 normalized scores of TD3+BC, AWAC, CQL, IQL from the reported scores in each paper (Fujimoto & Gu, 2021; Nair et al., 2020; Kumar et al., 2020b; Kostrikov et al., 2021a). 364

Figure 3 shows the learning curves of TD3+BC and the results verify the effectiveness of the proposed 365 method. After the pretraining period (indicated by the red vertical lines), one can notice that the 366 learning curves rapidly increase and achieve higher returns compared to the original methods. These 367 results suggest that our method accelerates training and enhances performance with only a few lines 368 of modifications on top of the baselines. Full graphs of TD3+BC are provided on Figure 12 in 369 Appendix G. 370

We also applied our method with different pretraining time-step ratios (e.g., 10% - 0.1M of 1M 371 steps) on TD3+BC over 5 seeds. The results are presented on Figure 4. Notably, regardless of the 372 pretraining time-step ratio, the proposed method demonstrates improved performance over different 373 pretraining rates. Overall, the pretraining time-step ratio of 3% yields a slightly higher total sum of 374 averaged scores while the results of the 10% ratio yield the lowest standard deviation. For all of the 375 other experiments in this paper, we use the pretraining time-step ratio of 10%. 376

Additional experiments are conducted on large-scale robotic manipulation tasks in Robomimic 377 benchmark, to verify the effectiveness of the proposed method for complex tasks. The proposed



Figure 5: Averaged success rate on the Robomimic benchmark. We evaluate both vanilla methods without pretraining (blue) and methods with pretraining (orange). 7 out of 8 cases depict notably improved performance in both environments.



Figure 6: **Averaged normalized scores in reduced datasets across data quality.** This figure shows the overall performance of our method across reduced dataset (*i.e.*, 1%, 3%, 10%, 30%, 100%) for three environments (*i.e.*, halfcheetah, hopper, walker2d) in D4RL. From the overall results, we conclude that our method guarantees better performance even in 10% of the datasets regardless of the data quality of the dataset, and even 1% for the *random* datasets and 3% for the *medium* datasets.

method is evaluated with tasks containing suboptimal transitions, where the proposed method improves the baselines on the D4RL benchmark. The averaged success rate of four offline RL baselines is reported in Figure 5 with and without applying the proposed method. As can be seen, all the methods with the proposed pretraining method are improved over the baselines in seven out of eight cases. Therefore, we conclude that the proposed method also effectively performs in solving more complex tasks. We also have conducted experiments on Adroit, 24-DOF environment, in Appendix D. The results also demonstrate that the proposed method is effective in solving complex tasks.

# 5.2 DATA EFFICIENCY ACROSS THE QUALITIES OF THE DATASETS

To validate the data efficiency of the proposed method, regardless of the dataset quality, we have examined the proposed method with TD3+BC in reduced datasets (*i.e.*, 1%, 3%, 10%, 30%, 100% of each dataset) across the data quality (*i.e.*, random, medium, medium replay, medium expert, expert) on D4RL over 5 seeds. To construct the reduced datasets, we have uniformly sampled the transition segments (*i.e.*, (s, a, r, s')) from each dataset. On the random datasets (a leftmost section in Figure 13), training with the proposed method with only 1% of the dataset outperforms the vanilla TD3+BC trained with full datasets at halfcheetah and warker2d environments. On the medium datasets (right to the random in Figure 13), the proposed method shows similar or improved results



Figure 7: Comparison of the proposed method with the other approaches on D4RL. We compare existing model-free offline RLs with our method to offline model-based RLs (i.e., MOPO, MOBILE) and a representation RL (*i.e.*, ACL) on D4RL over 3 seeds. The graph shows integrated results over *medium*, *medium-replay*, *medium-expert* datasets. The results show that the proposed method maintains the performance in reduced datasets, especially 1%, unlike the other approaches.



459 Figure 8: Average returns in reduced datasets across the dataset collection strategies. We 460 evaluate our method over different dataset collection strategies (i.e., SMM, RND, ICM). TD3 with our method outperforms the vanilla TD3 overall and even training with 10% of datasets outperforms the 461 vanilla TD3 with full datasets. From the results, we demonstrate that our method is data-efficient 462 regardless of the data distributions. 463

443

444

445

446

447

450

451

452

453

454

456

457

458

compared to the vanilla TD3+BC with full datasets by only using 3% of the datasets. On the other 466 datasets (i.e. medium-replay, medium-expert, and expert), the proposed method with 10% datasets 467 totally outperforms the vanilla TD3+BC with full datasets. From the overall results in Figure 6, we 468 conclude that our method guarantees better performance even in 10% of the datasets regardless of the 469 data quality of the dataset. 470

We also compare our method with offline model-based RL and representation approaches. We apply 471 our method to TD3+BC, AWAC, and CQL. We adopt MOPO (Yu et al., 2020) and MOBILE (Sun et al., 472 2023) as representatives of offline model-based RL, ACL (Yang & Nachum, 2021) as a representation 473 representative. We conduct the experiments on D4RL, medium, medium-replay, medium-expert 474 datasets over three seeds. Figure 7 shows integrated results over the datasets and Figure 14 shows 475 details. The results show that our method maintains the performance in reduced datasets compared 476 with the other approaches that spend extra training budget (e.g., training and forwarding the transition). 477 Especially in 1% datasets, CQL with our method largely outperforms the others. As a result, we 478 claim that our method is the most proper choice for data-efficient offline RL.

479 480

481

#### 5.3 DATA EFFICIENCY ACROSS THE DATA DISTRIBUTIONS

482 We assume that a small dataset would have a shifted distribution compared to a large dataset, for 483 instance, some small datasets have narrow support of visited states. Based on the assumption we have made, we evaluate our method across different dataset collection strategies since each dataset has a 484 different data distribution. In ExoRL (Yarats et al., 2022), we chose TD3 as a comparison algorithm 485 and SMM (Lee et al., 2019), RND (Burda et al., 2018), and ICM (Pathak et al., 2017), as walker



Figure 9: Effectiveness of the proposed method over narrow support of the visited states datasets. (Left) Visualized goal-reaching point mass agents and trajectories with different goals, portions, and 502 exploration methods. (Right) Averaged return of CQL trained with two datasets with and without the 503 proposed pretraining method.

501

486

walk task datasets. In Yarats et al. (2022), ICM shows the best performance, followed by RND, SMM, 507 and TD3 shows the best performance in ICM. We compare TD3 to TD3 with our method in reduced 508 datasets (*i.e.*, 1%, 10%, 100%) over three seeds. To construct reduced datasets, we select the data 509 from the front. Figure 8 shows the results. For all datasets, applying our method with only 10% of 510 datasets outperforms vanilla TD3 with full datasets. Especially in RND, even training with 1% of 511 datasets shows a significantly high average return. 512

Furthermore, we consider apoint mass maze environment in ExoRL to investigate whether our method 513 is effective even in narrow support of the visited states datasets. Figure 9 visualizes the trajectories of 514 each reduced dataset collected by DIAYN (Eysenbach et al., 2018), and Proto (Yarats et al., 2021a) 515 strategies (i.e., 1% of DIAYN, 7% of Proto). In comparison with Figure 2 in Yarats et al. (2022), 516 our reduced dataset settings cover more narrow support of visited states. The top right figure of 517 DIAYN shows that there are a few trajectories around the *top right goal* and the bottom left right 518 figure of Proto also shows that there are a few trajectories around the *bottom right goal* in Figure 9. 519 To demonstrate our method is effective even with a dataset with this shifted state distribution, we 520 evaluated the proposed method on reduced *point mass maze* datasets described in Figure 9 over 521 short (reach top right) and long (reach bottom right) goals with CQL. Figure 9 demonstrates that 522 our method shows significant performance even with narrow data distribution. From the results, we 523 conclude that our method is indeed more data-efficient than the other methods regardless of different choices of the data distribution. 524

525 526

#### 6 CONCLUSION

527 528

529 In this paper, we propose a simple yet effective data-efficient offline RL method that pretrains a 530 shared Q-network with the transition dynamics prediction task, maintaining reasonable performance 531 even with a small training dataset. To pretrain the Q-network, we design a novel shared network architecture that outputs predictions of the next state and Q-value. This structure makes our method 532 easy to apply to any existing offline RL algorithms and efficiently boosts data efficiency. 533

534 To demonstrate the effectiveness of the proposed strategy, we conduct experiments with various settings in offline RL. From the results, we demonstrate that our method significantly improves 536 the performance of existing offline RL algorithms over D4RL and Robomimic benchmarks. We also demonstrate that our method is indeed data-efficient across the different data qualities from D4RL and the different data distributions from ExoRL. We leave future work to expand our method 538 toward various offline RL problems, e.g., offline to online RL, goal-conditioned RL, and real-world applications.

# 540 REFERENCES

583

584

Rishabh Agarwal, Dale Schuurmans, and Mohammad Norouzi. An optimistic perspective on offline
reinforcement learning. In *International conference on machine learning*, pp. 104–114. PMLR, 2020.

Michael Ahn, Anthony Brohan, Noah Brown, Yevgen Chebotar, Omar Cortes, Byron David, Chelsea 546 Finn, Chuyuan Fu, Keerthana Gopalakrishnan, Karol Hausman, Alex Herzog, Daniel Ho, Jasmine Hsu, Julian Ibarz, Brian Ichter, Alex Irpan, Eric Jang, Rosario Jauregui Ruano, Kyle Jeffrey, Sally 547 Jesmonth, Nikhil J. Joshi, Ryan Julian, Dmitry Kalashnikov, Yuheng Kuang, Kuang-Huei Lee, 548 Sergey Levine, Yao Lu, Linda Luu, Carolina Parada, Peter Pastor, Jornell Quiambao, Kanishka Rao, 549 Jarek Rettinghouse, Diego Reyes, Pierre Sermanet, Nicolas Sievers, Clayton Tan, Alexander Toshev, 550 Vincent Vanhoucke, Fei Xia, Ted Xiao, Peng Xu, Sichun Xu, Mengyuan Yan, and Andy Zeng. Do 551 as i can, not as i say: Grounding language in robotic affordances. (arXiv:2204.01691), August 552 2022. doi: 10.48550/arXiv.2204.01691. URL http://arxiv.org/abs/2204.01691. 553 arXiv:2204.01691 [cs]. 554

557 Chethan Bhateja, Derek Guo, Dibya Ghosh, Anikait Singh, Manan Tomar, Quan Vuong, Yevgen Chebotar, Sergey Levine, and Aviral Kumar. Robotic offline rl from internet videos via value-function pre-training. (arXiv:2309.13041), September 2023. URL http://arxiv.org/abs/2309.13041. arXiv:2309.13041 [cs].

561 Anthony Brohan, Noah Brown, Justice Carbajal, Yevgen Chebotar, Xi Chen, Krzysztof Choro-562 manski, Tianli Ding, Danny Driess, Avinava Dubey, Chelsea Finn, Pete Florence, Chuyuan Fu, 563 Montse Gonzalez Arenas, Keerthana Gopalakrishnan, Kehang Han, Karol Hausman, Alexander Herzog, Jasmine Hsu, Brian Ichter, Alex Irpan, Nikhil Joshi, Ryan Julian, Dmitry Kalashnikov, 565 Yuheng Kuang, Isabel Leal, Lisa Lee, Tsang-Wei Edward Lee, Sergey Levine, Yao Lu, Henryk 566 Michalewski, Igor Mordatch, Karl Pertsch, Kanishka Rao, Krista Reymann, Michael Ryoo, Grecia Salazar, Pannag Sanketi, Pierre Sermanet, Jaspiar Singh, Anikait Singh, Radu Soricut, Huong 567 Tran, Vincent Vanhoucke, Quan Vuong, Ayzaan Wahid, Stefan Welker, Paul Wohlhart, Jialin 568 Wu, Fei Xia, Ted Xiao, Peng Xu, Sichun Xu, Tianhe Yu, and Brianna Zitkovich. Rt-2: Vision-569 language-action models transfer web knowledge to robotic control. (arXiv:2307.15818), July 570 2023a. doi: 10.48550/arXiv.2307.15818. URL http://arxiv.org/abs/2307.15818. 571 arXiv:2307.15818 [cs]. 572

- 573 Anthony Brohan, Noah Brown, Justice Carbajal, Yevgen Chebotar, Joseph Dabis, Chelsea Finn, 574 Keerthana Gopalakrishnan, Karol Hausman, Alex Herzog, Jasmine Hsu, Julian Ibarz, Brian Ichter, Alex Irpan, Tomas Jackson, Sally Jesmonth, Nikhil J. Joshi, Ryan Julian, Dmitry Kalashnikov, 575 Yuheng Kuang, Isabel Leal, Kuang-Huei Lee, Sergey Levine, Yao Lu, Utsav Malla, Deeksha 576 Manjunath, Igor Mordatch, Ofir Nachum, Carolina Parada, Jodilyn Peralta, Emily Perez, Karl 577 Pertsch, Jornell Quiambao, Kanishka Rao, Michael Ryoo, Grecia Salazar, Pannag Sanketi, Kevin 578 Sayed, Jaspiar Singh, Sumedh Sontakke, Austin Stone, Clayton Tan, Huong Tran, Vincent Van-579 houcke, Steve Vega, Quan Vuong, Fei Xia, Ted Xiao, Peng Xu, Sichun Xu, Tianhe Yu, and Brianna 580 Zitkovich. Rt-1: Robotics transformer for real-world control at scale. (arXiv:2212.06817), August 581 2023b. doi: 10.48550/arXiv.2212.06817. URL http://arxiv.org/abs/2212.06817. 582
  - Yuri Burda, Harrison Edwards, Amos Storkey, and Oleg Klimov. Exploration by random network distillation. *arXiv preprint arXiv:1810.12894*, 2018.
- Yevgen Chebotar, Quan Vuong, Karol Hausman, Fei Xia, Yao Lu, Alex Irpan, Aviral Kumar, Tianhe
  Yu, Alexander Herzog, Karl Pertsch, et al. Q-transformer: Scalable offline reinforcement learning
  via autoregressive q-functions. In *Conference on Robot Learning*, pp. 3909–3928. PMLR, 2023.
- Jie Cheng, Ruixi Qiao, Gang Xiong, Qinghai Miao, Yingwei Ma, Binhua Li, Yongbin Li, and Yisheng
   Lv. Scaling offline model-based rl via jointly-optimized world-action model pretraining. *arXiv* preprint arXiv:2410.00564, 2024.
- Marc Deisenroth and Carl E Rasmussen. Pilco: A model-based and data-efficient approach to policy
   search. In *Proceedings of the 28th International Conference on machine learning (ICML-11)*, pp. 465–472, 2011.

Philip J Ball, Laura Smith, Ilya Kostrikov, and Sergey Levine. Efficient online reinforcement learning
 with offline data. In *International Conference on Machine Learning*, pp. 1577–1594. PMLR, 2023.

594 595 596	Pierluca D'Oro, Max Schwarzer, Evgenii Nikishin, Pierre-Luc Bacon, Marc G Bellemare, and Aaron Courville. Sample-efficient reinforcement learning by breaking the replay ratio barrier. In <i>Deep Reinforcement Learning Workshop NeurIPS 2022</i> , 2022.
598 599	Benjamin Eysenbach, Abhishek Gupta, Julian Ibarz, and Sergey Levine. Diversity is all you need: Learning skills without a reward function. <i>arXiv preprint arXiv:1802.06070</i> , 2018.
600 601 602 603	Yunhai Feng, Nicklas Hansen, Ziyan Xiong, Chandramouli Rajagopalan, and Xiaolong Wang. Finetuning offline world models in the real world. October 2023. doi: 10.48550/arXiv.2310.16029. URL http://arxiv.org/abs/2310.16029.
604 605	Justin Fu, Aviral Kumar, Ofir Nachum, George Tucker, and Sergey Levine. D4rl: Datasets for deep data-driven reinforcement learning. <i>arXiv preprint arXiv:2004.07219</i> , 2020.
606 607 608 609 610	Scott Fujimoto and Shixiang (Shane) Gu. A minimalist approach to offline reinforcement learning. In <i>Advances in Neural Information Processing Systems</i> , volume 34, pp. 20132–20145. Curran Associates, Inc., 2021. URL https://proceedings.neurips.cc/paper_files/paper/2021/hash/a8166da05c5a094f7dc03724b41886e5-Abstract.html.
611 612	Scott Fujimoto, Herke Hoof, and David Meger. Addressing function approximation error in actor- critic methods. In <i>International conference on machine learning</i> , pp. 1587–1596. PMLR, 2018.
614 615 616 617	Scott Fujimoto, David Meger, and Doina Precup. Off-policy deep reinforcement learning with- out exploration. In <i>Proceedings of the 36th International Conference on Machine Learning</i> , pp. 2052–2062. PMLR, May 2019. URL https://proceedings.mlr.press/v97/ fujimoto19a.html.
618 619 620	Zhaohan Daniel Guo, Mohammad Gheshlaghi Azar, Bilal Piot, Bernardo A Pires, and Rémi Munos. Neural predictive belief representations. <i>arXiv preprint arXiv:1811.06407</i> , 2018.
621 622 623	Tuomas Haarnoja, Aurick Zhou, Kristian Hartikainen, George Tucker, Sehoon Ha, Jie Tan, Vikash Kumar, Henry Zhu, Abhishek Gupta, Pieter Abbeel, et al. Soft actor-critic algorithms and applications. <i>arXiv preprint arXiv:1812.05905</i> , 2018.
624 625 626 627	Danijar Hafner, Timothy Lillicrap, Jimmy Ba, and Mohammad Norouzi. Dream to control: Learning behaviors by latent imagination. (arXiv:1912.01603), 2019a. URL https://arxiv.org/abs/1912.01603.
628 629 630 631	Danijar Hafner, Timothy Lillicrap, Ian Fischer, Ruben Villegas, David Ha, Honglak Lee, and James Davidson. Learning latent dynamics for planning from pixels. In <i>International conference on machine learning</i> , pp. 2555–2565. PMLR, 2019b. URL https://proceedings.mlr.press/v97/hafner19a.html.
632 633 634	Danijar Hafner, Jurgis Pasukonis, Jimmy Ba, and Timothy Lillicrap. Mastering diverse domains through world models. <i>arXiv preprint arXiv:2301.04104</i> , 2023.
635 636 637	Nicklas Hansen, Xiaolong Wang, and Hao Su. Temporal difference learning for model predictive control. (arXiv:2203.04955), July 2022. doi: 10.48550/arXiv.2203.04955. URL http://arxiv.org/abs/2203.04955. arXiv:2203.04955 [cs].
639 640	Nicklas Hansen, Hao Su, and Xiaolong Wang. Td-mpc2: Scalable, robust world models for continuous control. <i>arXiv preprint arXiv:2310.16828</i> , 2023.
641 642 643 644	Kaiming He, Haoqi Fan, Yuxin Wu, Saining Xie, and Ross Girshick. Momentum contrast for unsupervised visual representation learning. In <i>Proceedings of the IEEE/CVF conference on computer vision and pattern recognition</i> , pp. 9729–9738, 2020.
645 646 647	Zhang-Wei Hong, Aviral Kumar, Sathwik Karnik, Abhishek Bhandwaldar, Akash Srivastava, Joni Pajarinen, Romain Laroche, Abhishek Gupta, and Pulkit Agrawal. Beyond uniform sampling: Of-fline reinforcement learning with imbalanced datasets. <i>Advances in Neural Information Processing Systems</i> , 36:4985–5009, 2023.

648 649 650	Ilya Kostrikov, Ashvin Nair, and Sergey Levine. Offline reinforcement learning with implicit q- learning. (arXiv:2110.06169), October 2021a. doi: 10.48550/arXiv.2110.06169. URL http: //arxiv.org/abs/2110.06169. arXiv:2110.06169 [cs].
652 653 654	Ilya Kostrikov, Denis Yarats, and Rob Fergus. Image augmentation is all you need: Regularizing deep reinforcement learning from pixels. (arXiv:2004.13649), March 2021b. URL http://arxiv.org/abs/2004.13649. arXiv:2004.13649 [cs, eess, stat].
655 656 657 658	Aviral Kumar, Justin Fu, Matthew Soh, George Tucker, and Sergey Levine. Stabilizing off-policy q- learning via bootstrapping error reduction. In Advances in Neural Information Processing Systems, volume 32. Curran Associates, Inc., 2019. URL https://proceedings.neurips.cc/ paper/2019/hash/c2073ffa77b5357a498057413bb09d3a-Abstract.html.
660 661	Aviral Kumar, Rishabh Agarwal, Dibya Ghosh, and Sergey Levine. Implicit under-parameterization inhibits data-efficient deep reinforcement learning. <i>arXiv preprint arXiv:2010.14498</i> , 2020a.
662 663 664 665	Aviral Kumar, Aurick Zhou, George Tucker, and Sergey Levine. Conservative q-learning for offline reinforcement learning. (arXiv:2006.04779), August 2020b. doi: 10.48550/arXiv.2006.04779. URL http://arxiv.org/abs/2006.04779. arXiv:2006.04779 [cs, stat].
666 667 668	Misha Laskin, Kimin Lee, Adam Stooke, Lerrel Pinto, Pieter Abbeel, and Aravind Srinivas. Rein- forcement learning with augmented data. <i>Advances in neural information processing systems</i> , 33: 19884–19895, 2020.
669 670 671	Lisa Lee, Benjamin Eysenbach, Emilio Parisotto, Eric Xing, Sergey Levine, and Ruslan Salakhutdinov. Efficient exploration via state marginal matching. <i>arXiv preprint arXiv:1906.05274</i> , 2019.
672 673 674	Seunghyun Lee, Younggyo Seo, Kimin Lee, Pieter Abbeel, and Jinwoo Shin. Offline-to-online reinforcement learning via balanced replay and pessimistic q-ensemble. In <i>Proceedings of the 5th Conference on Robot Learning</i> , pp. 1702–1712. PMLR, January 2022. URL https://
675	proceedings.mlr.press/v164/lee22d.html.
676 677 678 679	Sergey Levine, Aviral Kumar, George Tucker, and Justin Fu. Offline reinforcement learning: Tutorial, review, and perspectives on open problems. (arXiv:2005.01643), November 2020. doi: 10.48550/arXiv.2005.01643. URL http://arxiv.org/abs/2005.01643. arXiv:2005.01643 [cs, stat].
680 681 682 683 684	Ajay Mandlekar, Fabio Ramos, Byron Boots, Silvio Savarese, Li Fei-Fei, Animesh Garg, and Dieter Fox. Iris: Implicit reinforcement without interaction at scale for learning control from offline robot manipulation data. In 2020 IEEE International Conference on Robotics and Automation (ICRA), pp. 4414–4420. IEEE, 2020.
685 686 687	Ajay Mandlekar, Danfei Xu, Josiah Wong, Soroush Nasiriany, Chen Wang, Rohun Kulkarni, Li Fei- Fei, Silvio Savarese, Yuke Zhu, and Roberto Martín-Martín. What matters in learning from offline human demonstrations for robot manipulation. <i>arXiv preprint arXiv:2108.03298</i> , 2021.
688 689	Francisco S Melo and M Isabel Ribeiro. Q-learning with linear function approximation. 2007.
690 691 692	Ashvin Nair, Abhishek Gupta, Murtaza Dalal, and Sergey Levine. Awac: Accelerating online reinforcement learning with offline datasets. <i>arXiv preprint arXiv:2006.09359</i> , 2020.
693 694 695	Mitsuhiko Nakamoto, Simon Zhai, Anikait Singh, Max Sobol Mark, Yi Ma, Chelsea Finn, Aviral Kumar, and Sergey Levine. Cal-ql: Calibrated offline rl pre-training for efficient online fine-tuning. <i>Advances in Neural Information Processing Systems</i> , 36, 2024.
696 697 698	Abhishek Padalkar, Acorn Pooley, Ajinkya Jain, Alex Bewley, Alex Herzog, Alex Irpan, Alexander Khazatsky, Anant Rai, Anikait Singh, Anthony Brohan, et al. Open x-embodiment: Robotic learning datasets and rt-x models. <i>arXiv preprint arXiv:2310.08864</i> , 2023.
700 701	Deepak Pathak, Pulkit Agrawal, Alexei A Efros, and Trevor Darrell. Curiosity-driven exploration by self-supervised prediction. In <i>International conference on machine learning</i> , pp. 2778–2787. PMLR, 2017.

702 703 704	Rafael Rafailov, Kyle Beltran Hatch, Victor Kolev, John D Martin, Mariano Phielipp, and Chelsea Finn. Moto: Offline pre-training to online fine-tuning for model-based robot learning. In <i>Conference on Robot Learning</i> , pp. 3654–3671. PMLR, 2023.
705 706 707 708	Aravind Rajeswaran, Vikash Kumar, Abhishek Gupta, Giulia Vezzani, John Schulman, Emanuel Todorov, and Sergey Levine. Learning complex dexterous manipulation with deep reinforcement learning and demonstrations. <i>arXiv preprint arXiv:1709.10087</i> , 2017.
709 710 711	Max Schwarzer, Ankesh Anand, Rishab Goel, R Devon Hjelm, Aaron Courville, and Philip Bach- man. Data-efficient reinforcement learning with self-predictive representations. <i>arXiv preprint</i> <i>arXiv:2007.05929</i> , 2020.
712 713 714 715 716	Max Schwarzer, Nitarshan Rajkumar, Michael Noukhovitch, Ankesh Anand, Laurent Charlin, R De- von Hjelm, Philip Bachman, and Aaron C Courville. Pretraining representations for data-efficient reinforcement learning. <i>Advances in Neural Information Processing Systems</i> , 34:12686–12699, 2021.
717 718 719	Younggyo Seo, Kimin Lee, Stephen L. James, and Pieter Abbeel. Reinforcement learning with action- free pre-training from videos. In <i>International Conference on Machine Learning</i> , pp. 19561–19579. PMLR, 2022. URL https://proceedings.mlr.press/v162/seo22a.html.
720 721 722	Yihao Sun, Jiaji Zhang, Chengxing Jia, Haoxin Lin, Junyin Ye, and Yang Yu. Model-bellman inconsistency for model-based offline reinforcement learning. In <i>International Conference on Machine Learning</i> , pp. 33177–33194. PMLR, 2023.
723 724 725	Richard S Sutton. Dyna, an integrated architecture for learning, planning, and reacting. <i>ACM Sigart Bulletin</i> , 2(4):160–163, 1991.
726	Richard S Sutton, Andrew G Barto, et al. Introduction to reinforcement learning. vol. 135, 1998.
727 728 729 730	Octo Model Team, Dibya Ghosh, Homer Walke, Karl Pertsch, Kevin Black, Oier Mees, Sudeep Dasari, Joey Hejna, Tobias Kreiman, Charles Xu, et al. Octo: An open-source generalist robot policy. <i>arXiv preprint arXiv:2405.12213</i> , 2024.
731 732 733	Ziyu Wang, Tom Schaul, Matteo Hessel, Hado Hasselt, Marc Lanctot, and Nando Freitas. Dueling network architectures for deep reinforcement learning. In <i>International conference on machine learning</i> , pp. 1995–2003. PMLR, 2016.
734 735 736	Tengyang Xie, Nan Jiang, Huan Wang, Caiming Xiong, and Yu Bai. Policy finetuning: Bridg- ing sample-efficient offline and online reinforcement learning. <i>Advances in neural information</i> <i>processing systems</i> , 34:27395–27407, 2021.
737 738 739	Mengjiao Yang and Ofir Nachum. Representation matters: Offline pretraining for sequential decision making. In <i>International Conference on Machine Learning</i> , pp. 11784–11794. PMLR, 2021.
740 741 742	Rui Yang, Han Zhong, Jiawei Xu, Amy Zhang, Chongjie Zhang, Lei Han, and Tong Zhang. Towards robust offline reinforcement learning under diverse data corruption. <i>arXiv preprint</i> <i>arXiv:2310.12955</i> , 2023.
743 744 745 746	Denis Yarats, Rob Fergus, Alessandro Lazaric, and Lerrel Pinto. Reinforcement learning with prototypical representations. In <i>International Conference on Machine Learning</i> , pp. 11920–11931. PMLR, 2021a.
747 748 749	Denis Yarats, Ilya Kostrikov, and Rob Fergus. Image augmentation is all you need: Regularizing deep reinforcement learning from pixels. In <i>International conference on learning representations</i> , 2021b.
750 751 752 752	Denis Yarats, Amy Zhang, Ilya Kostrikov, Brandon Amos, Joelle Pineau, and Rob Fergus. Improving sample efficiency in model-free reinforcement learning from images. In <i>Proceedings of the aaai conference on artificial intelligence</i> , volume 35, pp. 10674–10681, 2021c.
753 754 755	Denis Yarats, David Brandfonbrener, Hao Liu, Michael Laskin, Pieter Abbeel, Alessandro Lazaric, and Lerrel Pinto. Don't change the algorithm, change the data: Exploratory data for offline reinforcement learning. <i>arXiv preprint arXiv:2201.13425</i> , 2022.

Tianhe Yu, Garrett Thomas, Lantao Yu, Stefano Ermon, James Y Zou, Sergey Levine, Chelsea Finn, and Tengyu Ma. Mopo: Model-based offline policy optimization. Advances in Neural Information Processing Systems, 33:14129–14142, 2020. Tianhe Yu, Aviral Kumar, Yevgen Chebotar, Karol Hausman, Chelsea Finn, and Sergey Levine. How to leverage unlabeled data in offline reinforcement learning. In International Conference on Machine Learning, pp. 25611–25635. PMLR, 2022. 

# A PRETRAINING Q-NETWORK FOR ONLINE RL (OFF-POLICY)

Inn	$\mathbf{ut}$ : Learning rate $\alpha$
Init	ialize parameters $\omega, \psi$ and a buffer $\mathcal{D}$
for	each gradient step <b>do</b>
	Uniformly sample a random action and collect a transition
	$a \sim U(a_{min}, a_{max})$
	$s' \sim p(s' s,a)$
	Update the buffer with a collected transition $\mathcal{D}_{i} = \mathcal{D}_{i} = \{(a_{i}, a_{j}, a_{j}, a_{j})\}$
	$D \leftarrow D \cup \{(s, a, r, s')\}$
	Sample a mini-batch $\mathcal{B} \sim \mathcal{D}$
	Compute the forward dynamics prediction error
	$\mathcal{L}_{pre}(arphi,\psi)=\sum_{a}(s'-(g_\psi\circ h_arphi)(s,a))^2$
	$(s,a,s') {\in} \mathcal{B}$
	Update weights of the shared network and forward network
	$(\alpha \leftarrow (\alpha - \alpha \nabla f (\alpha y)))  y \leftarrow y = \alpha \nabla f (\alpha y)$
	$r \sim r \sim \phi \sim pre(r, \tau),  \tau \sim \tau \sim \phi \sim pre(r, \tau)$
end	for
Out	tput: Pretrained weights $\varphi$ of the shared network, collected buffer $\mathcal{D}$
1	then 2 Destavising where for Outing DL (Off wallow) with one calledt a laterat
igor	itim 3 Pretraining phase for Online RL (Oπ-policy) with pre-collected dataset
Inp	<b>ut</b> : Dataset $\mathcal{D}_{pre}$ of transition $(s, a, s')$ , Learning rate $\alpha$
Init	ialize parameters $\varphi, \psi$
for	each gradient step <b>do</b>
	Sample a mini-batch $\mathcal{B} \sim \mathcal{D}_{pre}$
	Dating the loss tunction
	Denne tile toss function
	$\mathcal{L}_{nre}(\varphi,\psi) = \sum (s' - (q_{\psi} \circ h_{\omega})(s,a))^2$
	$\mathcal{L}_{pre}(\varphi,\psi) = \sum_{(s,a,s')\in\mathcal{B}} (s' - (g_{\psi} \circ h_{\varphi})(s,a))^2$
	$\mathcal{L}_{pre}(\varphi,\psi) = \sum_{(s,a,s')\in\mathcal{B}} (s' - (g_\psi \circ h_\varphi)(s,a))^2$ Take the evolution
	$\mathcal{L}_{pre}(\varphi,\psi) = \sum_{(s,a,s')\in\mathcal{B}} (s' - (g_\psi \circ h_\varphi)(s,a))^2$ Take the gradient descent step
	$\mathcal{L}_{pre}(\varphi, \psi) = \sum_{(s,a,s') \in \mathcal{B}} (s' - (g_{\psi} \circ h_{\varphi})(s,a))^2$ Take the gradient descent step $\varphi \leftarrow \varphi - \alpha \nabla_{\phi} \mathcal{L}_{pre}(\varphi, \psi),  \psi \leftarrow \psi - \alpha \nabla_{\psi} \mathcal{L}_{pre}(\varphi, \psi)$
end	$\mathcal{L}_{pre}(\varphi,\psi) = \sum_{(s,a,s')\in\mathcal{B}} (s' - (g_{\psi} \circ h_{\varphi})(s,a))^2$ Take the gradient descent step $\varphi \leftarrow \varphi - \alpha \nabla_{\phi} \mathcal{L}_{pre}(\varphi,\psi),  \psi \leftarrow \psi - \alpha \nabla_{\psi} \mathcal{L}_{pre}(\varphi,\psi)$ for

We extended our pretraining method to popular online off-policy RL methods by incorporating
the pretraining phase ahead of the main training phase. During the pretraining phase of the online
agent, a trajectory dataset was obtained by either initializing the replay buffer with actively collected
interaction data by uniformly sampling a random action or offline static dataset.

For experiments on online RL using an off-policy setting, we adopted soft actor-critic (SAC) Haarnoja
et al. (2018) and twin delayed deep deterministic policy gradient algorithm (TD3) Fujimoto et al.
(2018). We compare these algorithms with and without our pretraining method on OpenAI Gym
MuJoCo tasks. For a fair comparison, all algorithms were trained for 1 million time steps on each
task over 5 seeds.

Table 3 presents the results of the experiments following Algorithm 2 which collects the pretraining dataset by uniformly sampling random actions. Incorporating our pretraining phase shows better performance in more than half of the results. Additionally, we trained both SAC and TD3 with the pre-collected dataset from the D4RL for the pretraining phase along the Algorithm 3. Note that we only used the pre-collected dataset during the pretraining phase. Table 4 shows the best scores among the 5 datasets (i.e., random, medium, medium replay, medium expert, expert). Interestingly,

pretraining with the suboptimal-level dataset (medium-replay) shows better performance compared to the expert-level dataset.

Table 3: Results of Off-policy RL application on OpenAI gym MuJoCo tasks

	SAC	TD3
HalfCheetah-v2	$10065.77 \pm 621.80 \rightarrow 11005.51 \pm 374.14$	10644.63±190.42→11697.71±236.01
Hopper-v2	$3357.07 \pm 30.64 \rightarrow 1419.55 \pm 137.55$	$3365.08 \pm 94.69 \rightarrow 3454.83 \pm 129.34$
Walker2d-v2	$4279.67 \pm 509.51 \rightarrow 2697.92 \pm 674.29$	$4193.11 \pm 435.31 \rightarrow 4481.19 \pm 190.93$
Ant-v2	4191.17±986.11→4399.56 766.24	$5172.78 \pm 659.02 \rightarrow 4407.40 \pm 759.64$
Humanoid-v2	5545.70±85.00→479.09 83.86	5247.14±187.64→5816.16±199.25
Pusher-v2	-190.77±88.51→-133.96 29.00	$-22.94 \pm 0.52 \rightarrow -22.85 \pm 1.25$

Table 4: Results of Off-policy RL pretrain with the D4RL OpenAI gym MuJoCo datasets

	SAC	TD3
HalfCheetah-v2	10402.79±1675.67	11820.06±269.76
Hopper-v2	$3405.95 \pm 70.87$	$3465.25 \pm 149.87$
Walker2d-v2	$4785.15 \pm 247.37$	4559.38±1007.69

From the above experiments, we conjecture that pretrained online RL (off-policy) has limitations when they only exploit random action data for pretraining. A marginal state distribution induced by uniformly sampling random actions is close to the initial state distribution, limiting the diversity in the dataset and eventually leading to an increase in forward dynamics uncertainty. Consequently, there are fewer opportunities to learn the good features of forward dynamics with random action datasets than suboptimal-level datasets. This explains why Table 3 shows worse results than Table 4.

We also applied another approach introduced in section B to online RL settings. The results, shown in Table 5, indicate that more than half exhibit enhanced performance compared to reported scores in Table 3.

Table 5:	Results of	Off-policy R	RL with	Additional Loss

	SAC	TD3
HalfCheetah-v2	8498.68±3195.13	9588.53±866.30
Hopper-v2	3539.39±133.47	$3523.67 \pm 202.52$
Walker2d-v2	4847.86±135.52	$3819.68 \pm 552.84$
Ant-v2	3710.73±917.35	$5401.0 \pm 844.56$
Humanoid-v2	5576.98±106.31	$5489.73 \pm 38.28$
Pusher-v2	$-158.66 \pm 55.02$	$-25.47 \pm 34.00$

## **B** ANOTHER DESIGN CHOICE USING OUR SHARED Q-NETWORK STRUCTURE

In this section, we introduce another approach that also utilizes features of forward dynamics using the shared networks as in the previouse pretraining method. In this approach, we use the following modified loss that adds the forward model loss to the loss for the *Q*-function estimation:

$$\mathcal{L}_Q = \mathcal{L}_{TD} + \mathcal{L}_{dynamics} \tag{6}$$

In this way, the shared network is trained throughout the entire training period without the pretraining phase. We adopt TD3+BC for evaluation and the results are presented in table 6. On TD3+BC, this approach also outperforms almost all of the vanilla scores. Simply adding the supervised loss term of state prediction without any multiplier or technique demonstrates improved performance.
Consequently, we suggest that the proposed shared Q-network can be expanded in other directions and we expect that it holds significant potential for further research.

Table 6: Averaged normalized scores of TD3+BC with additional loss on D4RL benchmark. We
 depict increased scores compared to their original scores in blue color and report mean and standard
 deviations over 5 random seeds.

921							
922		Random	Medium	Medium Replay	Medium Expert	Expert	
923	HalfCheetah-v2	$11.45 \pm 0.51$	$48.23 \pm 0.33$	44.93±0.29	93.55±1.00	96.59±0.25	
924	Walker2d-v2	$13.46 \pm 6.58$	$70.86 \pm 2.17$ $82.65 \pm 1.65$	$90.39 \pm 7.34$ $86.11 \pm 1.54$	$113.44 \pm 0.35$ $111.88 \pm 0.63$	$113.28 \pm 0.20$ $110.98 \pm 0.22$	
925							
926							
927 (a) Ha	lfCheetah		(b)	Hopper		(c) Wal	ker2d
928							
929							
930							
931				1		Î	
932							
933							
934							
935			22				
936							
937							
938							
939 (d	l) Lift		(e	e) Can		(f) Walke	er Walk
940	0-1		and the second second				
941			Ó				
942		- M	THE ST				
943			AL POO				
944			1 C				
945	<b>V</b>						
940							
947			100 12 11				
949		7 A.					
950							
951							
952		Figure 10	). Illuctroi	tions of envi	ronmente		
953		i iguite It	). musu a		onnento.		
954							
955 C TASKS		SETS					
956 C TASKS	$\mathbf{D}$ AND $\mathbf{D}$ AIA	SEIS					
957 In this section	n we provid	e detailed	evnerime	ntal ceture f	or the tasks	and datace	te Illustrated
958 environments	s can be found	in Figure	10	mai setups I	or the tasks	and datase	is. musualeu
959	can be round						
960 C 1 D4PI							
961 C.I D4KL							
962 D4RL consis	sts of 8 separa	ate tasks.	In this w	ork, we utili	zed one of t	hem for the	e main experi-
963 ments; Open	AI Gym MuJo	Co contin	uous contr	ol tasks. It co	onsists of 4 d	ifferent envi	ronments (i.e.,
964 HalfCheetah,	Walker2d, Ho	opper, and	Ant) and 5	5 heterogeneo	ous datasets i	n terms of d	lata quality for

each environment. Each dataset is collected along the below strategies:

965

966

967 968

969

970

- Random (1M samples): Collected from a randomly initialized policy.
- Expert (1M samples): Collected from a policy trained to completion with SAC.
- Medium (1M samples): Collected from a policy trained to approximately 1/3 the performance of the expert.
  - Medium-Expert (almost 2M samples): A 50-50 split of medium and expert data.

• Medium-Replay (almost 3M samples): Collected from the replay buffer of a policy trained up to the performance of the medium agent.

All environments have the same episode limit of 1000 and the goal of each locomotion agent is to
run as fast as possible without falling to the ground. More detailed information can be found at
https://github.com/Farama-Foundation/D4RL.

C.2 ROBOMIMIC

Robomimic provides a large-scale and diverse collection of task demonstrations spanning multiple human or robotic demonstrations of varying quality. We considered machine-generated (MG) datasets generated by training an SAC agent for each task and then using intermediate policies to generate mixed-quality datasets. We selected this dataset for evaluation since our method demonstrated superior performance with suboptimal datasets on the D4RL benchmark. All environments have the same episode limit of 400. The goal of the Lift environment is lifting the cube above a certain height and the goal of the Can environment is placing the can into the corresponding container. More detailed information can be found at https://github.com/ARISE-Initiative/robomimic.

C.3 ExoRL

They provide exploratory datasets for 6 DeepMind Control Stuite domains (*i.e., Cartpole, Cheetah, Jaco Arm, Point Mass Maze, Quadruped, Walker*) and totally 19 tasks. For each domain, they
 collected datasets by running 9 unsupervised RL algorithms (*i.e., APS, APT, DIAYN, Disagreement, ICM, ProtoRL, Random, RND, SMM*) from URLB for total of 10M steps. More detailed information can be found at https://github.com/denisyarats/exorl?tab=readme-ov-file.

D EXPERIMENTS ON ADROIT IN D4RL

We conducted additional experiments on adroit in D4RL Fu et al. (2020) benchmark to validate that
our method can be adopted to different complex domains. An illustration of the Adroit environment
can be found in Figure 11. The Adroit domain involves controlling a 24-DoF robotic hand with
4 different control tasks (i.e., Pen, Door, Hammer, and Relocate) and 3 heterogeneous datasets as
following:

- Human: Collected with the 25 human demonstrations provided in the DAPG Rajeswaran et al. (2017) repository.
- Cloned: a 50-50 split between demonstration data and 2500 trajectories sampled from a behavioral cloned policy on the demonstrations. The demonstration trajectories are copied to match the number of behavioral cloned trajectories.
- Expert: Collected with 5000 trajectories sampled from an expert that solves the task, provided in the DAPG repository.



Figure 11: The tasks of Adroit. (top left) Pen - aligning a pen with a target orientation, (top right)
 Door - opening a door, (bottom left) Hammer - hammering a nail into a board, (bottom right)
 Relocate - moving a ball to a target position.

1026 For experiments, we compared AWAC, IQL, and TD3+BC with/without our pretraining method 1027 over 5 seeds. Table 7 yields averaged normalized scores for each task. Overall, learning with our 1028 pretraining phase demonstrates enhanced performance. From these results, we conclude that our 1029 method can be effective in complex domains not only tabletop but dexterous manipulation as well.

Table 7: Averaged normalized scores on Adroit. Left-hand side scores are scores of vanilla 1031 methods. Right-hand side scores are scores of baselines combined with our pretraining method. We 1032 depict increased scores compared to their original scores in blue color and report mean and standard 1033 deviations over 5 random seeds. 1034

		AWAC	IQL	TD3+BC
Human	Pen Hammer Door Relocate	$\begin{array}{c} 146.19{\pm}5.29{\rightarrow}157.60{\pm}5.28\\ 7.98{\pm}9.41{\rightarrow}36.95{\pm}35.13\\ 60.82{\pm}12.38{\rightarrow}29.96{\pm}22.43\\ 1.51{\pm}1.05{\rightarrow}3.91{\pm}2.21 \end{array}$	$\begin{array}{c} 101.87{\pm}14.34{\rightarrow}104.66{\pm}17.30\\ 14.33{\pm}5.22{\rightarrow}17.78{\pm}9.27\\ 6.74{\pm}1.31{\rightarrow}5.81{\pm}3.20\\ 1.20{\pm}1.05{\rightarrow}1.52{\pm}1.11 \end{array}$	$\begin{array}{c} 20.32{\pm}5.97{\rightarrow}20.78{\pm}10.93\\ 2.40{\pm}0.16{\rightarrow}2.38{\pm}0.17\\ -0.09{\pm}0.00{\rightarrow}{-}0.04{\pm}0.04\\ -0.29{\pm}0.01{\rightarrow}{-}0.18{\pm}0.13 \end{array}$
Cloned	Pen Hammer Door Relocate	$145.37 \pm 4.19 \rightarrow 144.48 \pm 3.42 \\ 10.37 \pm 7.88 \rightarrow 12.61 \pm 8.66 \\ 2.95 \pm 2.97 \rightarrow 9.59 \pm 7.73 \\ 0.04 \pm 0.09 \rightarrow 0.18 \pm 0.21 \\ \end{array}$	$\begin{array}{c} 98.38{\pm}16.13{\rightarrow}97.76{\pm}16.90\\ 8.94{\pm}2.07{\rightarrow}11.38{\pm}4.46\\ 5.61{\pm}3.02{\rightarrow}5.00{\pm}1.44\\ 0.91{\pm}0.45{\rightarrow}1.06{\pm}0.40 \end{array}$	$\begin{array}{c} 39.69{\pm}18.95{\rightarrow}48.18{\pm}11.27\\ 0.59{\pm}0.17{\rightarrow}1.17{\pm}0.61\\ -0.23{\pm}0.11{\rightarrow}-0.03{\pm}0.03\\ -0.02{\pm}0.09{\rightarrow}-0.13{\pm}0.09 \end{array}$
Expert	Pen Hammer Door Relocate	$\begin{array}{c} 163.99 {\pm}1.19 {\rightarrow}163.73 {\pm}1.88 \\ 130.08 {\pm}1.30 {\rightarrow}130.04 {\pm}0.48 \\ 106.67 {\pm}0.28 {\rightarrow}106.95 {\pm}0.16 \\ 109.70 {\pm}1.32 {\rightarrow}111.27 {\pm}0.35 \end{array}$	$\begin{array}{c} 148.38 {\pm} 2.46 {\rightarrow} 147.79 {\pm} 3.06 \\ 129.46 {\pm} 0.42 {\rightarrow} 129.50 {\pm} 0.36 \\ 106.45 {\pm} 0.29 {\rightarrow} 106.71 {\pm} 0.28 \\ 110.13 {\pm} 1.52 {\rightarrow} 109.82 {\pm} 1.45 \end{array}$	$\begin{array}{c} 131.73 {\pm} 19.15 {\rightarrow} 141.10 {\pm} 10.28\\ 33.36 {\pm} 34.61 {\rightarrow} 59.76 {\pm} 52.35\\ 0.99 {\pm} 0.83 {\rightarrow} 0.87 {\pm} 1.48\\ 0.57 {\pm} 0.33 {\rightarrow} 0.22 {\pm} 0.13 \end{array}$
Total		885.67±47.35→907.26±87.94	$732.40 \pm 48.27 \rightarrow 738.79 \pm 59.23$	$229.03 \pm 80.40 \rightarrow 274.08 \pm 87.49$

#### E IMPLEMENTATION DETAILS

In this section, we provide detailed implementation setups for extensive experiments. Since we 1051 suggest a plug-and-play pretraining method for popular offline RL methods, we reuse open-source 1052 code for comparative results:  $TD3+BC^1$ ,  $IQL^2$ , AWAC<sup>3</sup>, and  $CQL^4$  for D4RL. We use off-the-shelf 1053 offline methods in the official repository<sup>5</sup> for the Robomimic environment. We only use open-source 1054 baselines which use PyTorch for fair comparisons. On the D4RL, we train each agent with 1M 1055 gradient steps for each environment over 5 seeds. Also, we evaluate each agent with 5 rollouts every 1056 5k gradient steps for TD3+BC, AWAC, and CQL and 10k gradient steps for IQL. We report the best 1057 scores for all tables and figures. On the Robomimic, we train each agent with 200k gradient steps for 1058 each environment over 5 seeds. Also, we evaluate each agent with 50 rollouts over 5 seeds. For all 1059 experiments, we used RTX-A5000 GPU for training and evaluation.

F DISCUSSIONS

1030

1047 1048 1049

1050

1061

1062

1063 In this section, we address the potential concerns regarding our method's novelty since it closely 1064 connects with prior approaches in relevant fields. We provide our detailed discussions in separate subsections of each topic. 1066

**Representation Learning.** Over recent years, the field has observed a significant amount of literature 1067 working on predictive representation in RL. Concerning the similarity with prior works, we claim that 1068 the idea of pretraining shared Q-network for improving data efficiency is remarkable. Our method 1069 pretrains the neural networks with the next state prediction objective to improve an underlying RL 1070 agent's performance and data efficiency similar to (Schwarzer et al., 2020; Guo et al., 2018). However, 1071 Schwarzer et al. (2020) has proposed an online training method in a self-supervised learning manner 1072 whereas our method considers supervised learning for pretraining. Since the self-predictive task in Schwarzer et al. (2020) is conducted in latent space, representation learning is essentially involved 1074 with the task. 1075

<sup>1076</sup> <sup>1</sup>https://github.com/sfujim/TD3\_BC

<sup>&</sup>lt;sup>2</sup>https://github.com/Manchery/iql-pytorch 1077

<sup>&</sup>lt;sup>3</sup>https://github.com/hari-sikchi/AWAC 1078

<sup>&</sup>lt;sup>4</sup>https://github.com/young-geng/CQL 1079

<sup>&</sup>lt;sup>5</sup>https://github.com/ARISE-Initiative/robomimic

Therefore, adopting advanced training techniques including data augmentation (Yarats et al., 2021b)
and the use of a target encoder (He et al., 2020) significantly affect the RL agent's performance.
Additionally, Schwarzer et al. (2020) suggests a self-supervised representation learning with the
latent transition prediction task in the online RL regime. In comparison, our method alleviates
an introduction of extra techniques other than the shared network architecture, proving superior
performance in offline RL benchmarks of diverse environments, e.g. locomotion and manipulation
tasks.

1087 Guo et al. (2018) has presented an unsupervised learning method that encodes the *belief state* capturing 1088 sufficient information of the hidden true state from a past interaction history. In other words, the main 1089 interest of Guo et al. (2018) is how the neural network architecture trained with unsupervised learning 1090 extracts adequate information concerning the true state in POMDP, not how the underlying RL method given rich representation performs decision-making problem well. Specifically, the network 1091 architecture in Guo et al. (2018) is based on GRU, RNN based sequential network, and predicts a 1092 next observation  $o_{t+1}$  using action  $a_t$  and a belief state  $b_t$  that contains the partial information of 1093 the previous trajectory. Conversely, our method is implemented on MLP with the shared network 1094 architecture and predicts the next state  $s_{t+1}$  using current state  $s_t$  and action  $a_t$  without a past history. 1095

1096 Model-based RL. One might argue that our method lacks novelty with the idea of training a 1097 neural network with the transition dynamics prediction task. Obviously, the idea of approximating the transition dynamics (Sutton, 1991) for downstream RL training is not what we first suggest. 1098 However, we contend that our method has a few refuting viewpoints with previous similar works. 1099 TDMPC (Hansen et al., 2022) and TDMPC2 (Hansen et al., 2023) are model-based single and 1100 multi-task RL approaches, which recursively feed the output of the same network (i.e. the encoder 1101 and task embedding network) for the transition model and value learning. The outputs of the shared 1102 backbone networks correspond to the latent representation and task embedding vector, respectively, 1103 and most latent model-based RL approaches including TDMPC reuse the outputs for the transition 1104 model and value learning. On the other hand, our method presents a shared network architecture 1105 resembling the dueling architecture (Wang et al., 2016) to pretrain the shared backbone network with 1106 a separated stream (a header) of the transition model and Q-network. Additionally, this paper presents 1107 a two-phase training scheme: the transition model combined with the shared network is trained with the transition dynamics prediction task in the first phase and the Q-network, consisting of an MLP 1108 1109 header and the shared network initialized with the parameter of the shared network in the first phase, is trained with the downstream RL value learning task in the second phase. 1110

1111 JOWA (Cheng et al., 2024) is an offline world model for multi-task RL with a shared Transformer 1112 backbone network for sequential a next-token prediction task. By modeling the decision-making 1113 problem to the sequential token prediction task, the backbone network, tokenizer, and header are 1114 trained in a supervised manner with the offline dataset. While the main purpose of JOWA is scaling an offline world model across multiple tasks with generalized performance over unseen tasks, this 1115 paper intends to improve the data efficiency of conventional offline RL approaches in single-task 1116 RL. Furthermore, our method alleviates additional training after offline RL training with a novel 1117 two-phase training strategy while JOWA allows few-shot fine-tuning for sample efficient transfer 1118 with a multi-game environment. Even with a similar purpose of data efficiency, our method entails a 1119 minimal algorithmic change with a consistent training budget compared to previous approaches. 1120

Dreamer (Hafner et al., 2023) has brought a notable advancement in model-based RL. Dreamer 1121 suggests a world model for decision-making with a considerate design of the latent transition model 1122 and reconstructive objective. Since jointly learning an accurate world model and actor in a multi-task 1123 environment is challenging, the expensive cost of collecting samples often becomes problematic. In 1124 contrast, our method does not necessitate extra modifications of conventional offline RL and proves 1125 its sufficient performance gains in comprehensive experiments. Considering previous improvements 1126 in representation learning usually involve state-of-the-art design choices (e.g. data augmentation), 1127 this paper would contribute to reasonable architectural achievements for researchers by presenting a 1128 minimal training structure with verified performance profit.

- 1129
- 1130

1131

#### 1134 LEARNING CURVES G 1135

1170

1171



1164 In this section, We pretrained TD3+BC and froze it except for the last linear layer during the remaining 1165 learning time. The blue-colored scores indicate improved scores from the reported scores from the 1166 original TD3+BC. Although only the last linear layer of the pretrained TD3+BC was trained and the 1167 shared network was frozen, it shows better performance than the vanilla CQL. Moreover, it shows 1168 better performance than the others over the suboptimal level of the datasets (i.e., random, medium, medium replay). 1169

Table 8: Results of pretrained TD3+BC which approximated with linear Q function.

1172							
1173			AWAC	CQL	IQL	TD3+BC	freezed TD3+BC
1174		HalfCheetah	2.2	$21.7 {\pm} 0.9$		$10.2 \pm 1.3$	$6.03 \pm 2.65$
1175	Random	Hopper Walker2d	9.6 5.1	$10.7 \pm 0.1$ 2.7 ± 1.2		$11.0\pm0.1$ $1.4\pm1.6$	$11.59 \pm 10.56$ $7.18 \pm 0.58$
1176		HalfCheetah	37.4	37.2±0.3	47.4	42.8±0.3	42.64±1.19
1177	Medium	Hopper	72.0	$44.2 \pm 10.8$	66.4	$99.5 \pm 1.0$	$67.16 \pm 3.56$
1178		Walker2d	30.1	$57.5 \pm 8.3$	78.3	$79.7 \pm 1.8$	$72.03 \pm 0.78$
1179	Medium Replay	HalfCheetah Hopper		$41.9 \pm 1.1$ 28.6 ± 0.9	44.2 94.7	$43.3 \pm 0.5$ 31.4 \pm 3.0	$40.21 \pm 0.79$
1180	Wedium Replay	Walker2d		$15.8 \pm 2.6$	73.9	$25.2\pm5.1$	$41.02 \pm 12.05$
1181		HalfCheetah	36.8	27.1±3.9	86.7	97.9±4.4	47.35±8.73
1182	Medium Expert	Hopper	80.9	$111.4 \pm 1.2$	91.5	$112.2 \pm 0.2$	$95.07 \pm 15.27$
1183		Walker2d	42.7	68.1±13.1	109.6	$101.1 \pm 9.3$	74.75±0.59
1107	_	HalfCheetah	78.5	$82.4 \pm 7.4$		$105.7 \pm 1.9$	$61.93 \pm 10.71$
1104	Expert	Hopper	85.2	$111.2\pm 2.1$		$112.2 \pm 0.2$	$113.13 \pm 0.39$
1185		Walker2d	57.0	$103.8 \pm 7.6$		$105.7 \pm 2.7$	57.14±44.96
1186	Total			764.3±61.5		979.3±33.4	$801.64 \pm 132.34$
1187							

In this section, we provide the full results of learning curves in the section 5.1 for further information.

#### 1188 EXPERIMENTS WITH VARIOUS AMOUNT OF DATA Ι 1189

In this section, we provide more details in section 5.2 of the dataset size. We conducted each experiment with the same settings in subsection 5.1 over 5 seeds and reported the results that exhibit averaged normalized scores.



Figure 13: Averaged normalized scores across dataset optimal quality and sizes. This figure compares the performance of our method with TD3+BC in reduced datasets (*i.e.*, 1%, 3%, 10%, 30%, 100% of each dataset) to vanilla TD3+BC across the data quality (i.e., random, medium, medium replay, medium expert, expert) on D4RL. From the overall results (Bottom Right), we conclude that 1215 our method guarantees better performance even in 10% of the datasets regardless of the data quality of the dataset.



1231 Figure 14: Comparison with offline model-based RL and representation approaches. We 1232 compare TD3+BC, AWAC, CQL with ours to offline model-based RLs (*i.e.*, MOPO, Mobile) and 1233 a representation RL (i.e., ACL) on D4RL over 3 seeds. The gragh shows results over medium, medium-replay, medium-expert datasets. The results show that our method maintains the performance 1234 in reduced datasets, especially 1%, unlike the other approaches. 1235

1236

1190

1191

1192

1212

1213

1214

1216

- 1237
- 1238
- 1239
- 1240
- 1241

		w/o pretrain	w/ pretrain 10%	w/ pretrain 30%	w/ pretrain
	HalfCheetah	2.2	$9.71 \pm 3.08$	36.37±1.47	51.10±0.89
Random	Hopper	9.6	$97.05 \pm 3.24$	$93.35 \pm 6.32$	59.47±33.79
	Walker2d	5.1	$8.57 \pm 0.47$	$8.36 \pm 1.30$	$13.11 \pm 3.91$
	HalfCheetah	37.4	55.47±1.52	56.64±2.68	54.63±1.45
Medium	Hopper	72.0	$101.28 \pm 0.78$	$101.32 \pm 0.20$	$101.73 \pm 0.20$
	Walker2d	30.1	$95.14 \pm 1.46$	$91.38 \pm 1.37$	89.51±0.88
	HalfCheetah		$51.00 \pm 0.69$	$52.12 \pm 0.76$	55.75±1.30
Medium Replay	Hopper		$103.67 \pm 1.81$	$107.69 \pm 1.71$	$106.67 \pm 0.59$
	Walker2d		$104.10 \pm 1.57$	$105.42 \pm 1.97$	$100.31 \pm 2.11$
	HalfCheetah	36.8	83.18±1.69	86.55±0.94	90.05±1.89
Medium Expert	Hopper	80.9	$113.01 \pm 0.71$	$113.34 \pm 0.09$	$113.23 \pm 0.22$
	Walker2d	42.7	$117.26 \pm 1.77$	$114.68 {\pm} 2.18$	$111.88 \pm 0.28$
	HalfCheetah	78.5	91.54±1.04	93.46±0.54	93.48±0.11
Expert	Hopper	85.2	$113.02 \pm 0.17$	$113.18 \pm 0.20$	$112.86 \pm 0.10$
-	Walker2d	57.0	$117.92 \pm 2.07$	$112.55 {\pm} 0.56$	$111.22 \pm 0.35$
Total			$1261.90 \pm 22.05$	$1286.43 \pm 22.28$	$1265.01 \pm 48.0$

# Table 9: Results of pretrained AWAC over various size.

Table 10: Results of pretrained IQL over varying dataset sizes.

		w/o pretrain	w/ pretrain 10%	w/ pretrain 30%	w/ pretrain
Random	HalfCheetah Hopper Walker2d		$6.92 \pm 0.63$ $8.17 \pm 0.54$ $8.26 \pm 0.64$	$12.65 \pm 2.53$ $9.93 \pm 1.19$ $9.08 \pm 0.96$	$\substack{18.28 \pm 1.02 \\ 10.67 \pm 0.41 \\ 8.88 \pm 0.71 }$
Medium	HalfCheetah Hopper Walker2d	47.4 66.4 78.3	$46.51 \pm 0.18$ $75.72 \pm 3.23$ $82.62 \pm 1.03$	$47.87 \pm 0.21$ $80.76 \pm 3.51$ $83.89 \pm 1.69$	$\substack{48.85 \pm 0.16 \\ 78.62 \pm 2.21 \\ 83.63 \pm 1.14}$
Medium Replay	HalfCheetah Hopper Walker2d	44.2 94.7 73.9	$33.49 \pm 1.26$ $80.59 \pm 8.25$ $39.08 \pm 10.42$	41.16±0.50 91.08±3.67 75.33±4.17	45.48±0.17 99.43±1.71 87.95±1.68
Medium Expert	HalfCheetah Hopper Walker2d	86.7 91.5 109.6	87.44±2.52 93.89±10.67 111.23±0.83	$93.66 \pm 0.46$ $91.05 \pm 18.78$ $111.65 \pm 0.93$	95.25±0.14 105.77±11.31 112.09±0.93
Expert	HalfCheetah Hopper Walker2d		$77.85 \pm 3.82$ $109.16 \pm 3.25$ $113.76 \pm 2.55$	95.88±0.44 112.85±1.30 112.53±1.35	97.40±0.13 113.34±0.46 112.80±1.08
Total			974.68±49.84	1069.36±41.69	1118.46±23.25