

# CONTRASTIVE VALUE LEARNING: IMPLICIT MODELS FOR SIMPLE OFFLINE RL

**Anonymous authors**

Paper under double-blind review

## ABSTRACT

Model-based reinforcement learning (RL) methods are appealing in the offline setting because they allow an agent to reason about the consequences of actions without interacting with the environment. Prior methods learn a 1-step dynamics model, which predicts the next state given the current state and action. These models do not immediately tell the agent which actions to take, but must be integrated into a larger RL framework. Can we model the environment dynamics in a different way, such that the learned model does directly indicate the value of each action? In this paper, we propose Contrastive Value Learning (CVL), which learns an implicit, multi-step model of the environment dynamics. This model can be learned without access to reward functions, but nonetheless can be used to directly estimate the value of each action, without requiring any TD learning. Because this model represents the multi-step transitions implicitly, it avoids having to predict high-dimensional observations and thus scales to high-dimensional tasks. Our experiments demonstrate that CVL outperforms prior offline RL methods on complex continuous control benchmarks.

## 1 INTRODUCTION

While the offline RL setting is relevant to many real-world applications where the ability for online data collection is limited, it often requires RL algorithms to find policies that are not well-supported by the training data. Instead of learning via trial-and-error, offline RL algorithms must leverage logged historical data to learn about the outcome of different actions, potentially by capturing environment dynamics as a proxy signal. Many prior approaches for this offline RL setting have been proposed, whether in model-free (Wu et al., 2019; Fujimoto et al., 2019; Kumar et al., 2020) or model-based (Kidambi et al., 2020; Yu et al., 2021) settings. Our focus will be on those that address this prediction problem head-on: by learning a predictive model of the environment which can be used in conjunction with most model-free algorithms.

Prior model-based methods (Yu et al., 2020b; Argenson and Dulac-Arnold, 2020; Kidambi et al., 2020; Yu et al., 2021) learn a model that predicts the observation at the next time step. This model is then used to generate synthetic data that can be passed to an off-the-shelf RL algorithm. While these approaches can work well on some benchmarks, they can be complex and expensive: the model must predict high-dimensional observations, and determining the value of an action may require unrolling the model for many steps. Learning a model of the environment has not made the RL problem any simpler. Moreover, as we will show later in the paper, the environment dynamics are intertwined with the policy inside the value function; model-based methods aim to decouple these quantities by separately estimating them. On the other hand, we show that one can directly learn a long-horizon transition model for a given policy, which is then used to estimate the value function. A natural use case for learning an occupancy measure from unlabelled data is multi-task pretraining, where the implicit dynamics model is trained on trajectory data across a collection of tasks, often exhibiting positive transfer properties. As we demonstrate in our experiments, this multi-task occupancy measure can then be finetuned using reward-labelled states on the task of interest, greatly improving performance upon existing pretraining methods as well as *tabula rasa* approaches.

In this paper, we propose to learn a different type of model for offline RL, a model which (1) will not require predicting high-dimensional observations and (2) can be directly used to estimate Q-values without requiring either model-based rollouts or model-free temporal difference learning. Precisely,

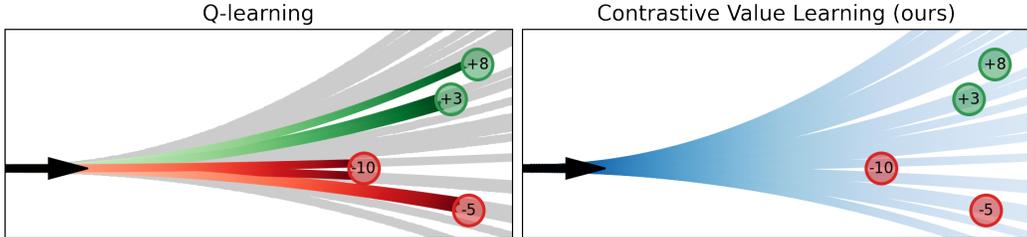


Figure 1: **Contrastive Value Learning**: A stylized illustration of trajectories (grey) and the rewards at future states (e.g., +8, -5). (Left) Q-learning estimates Q-values by “backing up” the rewards at future states. (Right) Our method learns the Q-values by fitting an implicit model to estimate the likelihoods of future states (blue), and taking the reward-weighted average of these likelihoods.

we will learn an implicit model of the discounted state occupancy measure, which answers the question “where will the agent be in the *time-averaged* future?” We will learn this implicit model via contrastive learning, treating it as a classifier rather than a generative model of observations. Once learned, we predict the likelihood of reaching every reward-labeled state. By weighting these predictions by the corresponding rewards, we form an unbiased estimate of the Q-function. Whereas methods like Q-learning estimate the Q-function of a state “backing up” reward values, our approach goes in the opposite direction, “propagating forward” predictions about where the agent will go.

We name our proposed algorithm Contrastive Value Learning (CVL). CVL is a simple algorithm for offline RL which learns the future state occupancy measure using contrastive learning and re-weights it with the future reward samples to construct a quantity proportional to the true value function. Because CVL represents multi-step transitions implicitly, it avoids having to predict high-dimensional observations and thus scales to high-dimensional tasks. Using the same algorithm, we can handle settings where reward-free data is provided, which cannot be directly handled by classical offline RL methods such as FQI (Munos, 2003) or BCQ (Fujimoto et al., 2019). We compare our proposed method to competitive offline RL baselines, notably CQL (Kumar et al., 2020) and CQL+UDS (Yu et al., 2022) on an offline version of the multi-task Metaworld benchmark (Yu et al., 2020a), and find that CVL greatly outperforms the baseline approaches as measured by the `reliable` library (Agarwal et al., 2021b). Additional experiments on image-based tasks from this same benchmark show that our approach scales to high-dimension tasks more seamlessly than the baselines. We also conduct a series of ablation experiments highlighting critical components of our method.

## 2 RELATED WORKS

Prior work has given rise to multiple offline RL algorithms, which often rely on behavior regularization in order to be well-supported by the training data. The key idea of offline RL methods is to balance interpolation and extrapolation errors, while ensuring proper diversity of out-of-dataset actions. Popular offline RL algorithms such as BCQ and CQL rely on a behavior regularization loss (Wu et al., 2019) as a way to control the extrapolation error. This regularization term ensures that the learned policy is well-supported by the data, i.e. does not stray too far away from the logging policy. The major issue with current offline RL algorithms is that they fail to fully capture the entire distribution over state-action pairs present in the training data.

To directly learn a value function using policy or value iteration, one needs to have information about the transition model in the form of sequences of state-action pairs, as well as the reward emitted by this transition. However, in some real-world scenarios, the reward might only be available for a small subset of data. For instance, in the case of recommending products available in an online catalog to the user, the true long-term reward (user buys the product) is only available for users who have browsed the item list for long enough and have purchased a given item. It is possible to decompose the value function into reward-dependent and reward-free parts, as was done by (Barreto et al., 2016) through the successor representation framework (Dayan, 1993). More recent approaches (Janner et al., 2020; Eysenbach et al., 2020; 2022) use a generative model to learn the occupancy measure over future states for each state-action pair in the dataset; its expectation corresponds to the successor representation. However, learning an explicit multi-step model such as (Janner et al., 2020) can be unstable due to the bootstrapping term in the temporal difference loss. Similarly to model-based

approaches, our method will learn a reward-free representation of the world, but will do so without having to predict high-dimensional observations and without having to do costly autoregressive rollouts. Thus, while our critic is trained without requiring rewards, it is much more similar to a value function than a standard 1-step model.

Learning a conditional probability distribution over a highly complex space can be a challenging task, which is why it is often easier to instead approximate it using a density ratio specified by an inner product in a much lower-dimensional latent space. To learn an occupancy measure over future states without passing via the temporal difference route, one can use noise-contrastive estimation (NCE, [Gutmann and Hyvärinen, 2010](#); [Oord et al., 2018](#)) to approximate the corresponding log ratio of densities as an implicit function. Contrastive learning was originally proposed as an alternative to classical maximum likelihood estimation, but has since then seen successes in static self-supervised learning ([He et al., 2020](#); [Chen et al., 2020](#)). In reinforcement learning, NCE was shown to improve the robustness of state representations to exogenous noise ([Srinivas et al., 2020](#); [Mazouze et al., 2020](#); [Agarwal et al., 2021a](#)) and, more recently, to be an efficient replacement for traditional goal-conditioned methods ([Eysenbach et al., 2022](#)).

### 3 PRELIMINARIES

**Reinforcement learning** We assume a Markov decision process  $M$  defined by the tuple  $M = \langle \mathcal{S}, S_0, \mathcal{A}, \mathcal{T}, r, \gamma \rangle$ , where  $\mathcal{S}$  is a state space,  $S_0 \subseteq \mathcal{S}$  is the set of starting states,  $\mathcal{A}$  is an action space,  $\mathcal{T} = \mathbb{P}[\cdot | s_t, a_t] : \mathcal{S} \times \mathcal{A} \rightarrow \Delta(\mathcal{S})$  is a one-step transition function<sup>1</sup>,  $r : \mathcal{S} \times \mathcal{A} \rightarrow [r_{\min}, r_{\max}]$  is a reward function and  $\gamma \in [0, 1)$  is a discount factor. The system starts in one of the initial states  $s_0 \in S_0$ . At every timestep  $t = 1, 2, 3, \dots$ , the policy  $\pi : \mathcal{S} \rightarrow \Delta(\mathcal{A})$ , samples an action  $a_t \sim \pi(\cdot | o_t)$ . The environment transitions into a next state  $s_{t+1} \sim \mathcal{T}(\cdot | s_t, a_t)$  and emits a reward  $r_t = r(s_t, a_t)$ . The aim is to learn a Markovian policy  $\pi(a | s)$  that maximizes the discounted sum of returns over an episode of length  $H$ :

$$\max_{\pi \in \Pi} \mathbb{E}_{\mathbb{P}_{0:H}^{\pi}, S_0} \left[ \sum_{t=0}^H \gamma^t r(s_t, a_t) \right], \quad (1)$$

where  $\mathbb{P}_{t:t+K}^{\pi}$  denotes the joint distribution of  $\{s_{t+k}, a_{t+k}\}_{k=1}^K$  obtained by executing  $\pi$  in the environment for  $K$  timesteps starting at timestep  $t$ .

We assume that the offline RL algorithm cannot interact with the environment, but instead must learn from an offline dataset of logged trajectories  $\{(s_0, a_0, s_1, a_1, \dots)\}$ . Value-based RL algorithms maximize cumulative episodic rewards by estimating the state-action value function under a policy  $\pi$ :

$$Q^{\pi}(s_t, a_t) = \mathbb{E}_{\mathbb{P}_t^{\pi}} \left[ \sum_{k=1}^H \gamma^k r(s_{t+k}, a_{t+k}) | s_t, a_t \right], \quad (2)$$

for  $s_t \in \mathcal{S}, a_t \in \mathcal{A}$ . Alternatively, the value function can be written as the expectation of the reward over the discounted occupancy measure:

$$Q^{\pi}(s_t, a_t) = \frac{1}{1 - \gamma} \mathbb{E}_{s, a \sim \mathbb{P}_{t:H}^{\pi}(s_t, a_t), \pi(s)} [r(s, a)], \quad (3)$$

where  $\mathbb{P}_{t:H}^{\pi}(s | s_t, a_t) = (1 - \gamma) \sum_{\Delta t=1}^H \gamma^{\Delta t-1} \mathbb{P}[S_{t+\Delta t} = s | s_t, a_t; \pi]$  as defined in [Janner et al. \(2020\)](#). Note that the occupancy measure can equivalently be re-written in terms of the geometric distribution over the time interval  $[0, \infty)$  for infinite-horizon rollouts:

$$\mathbb{P}_{0:\infty}^{\pi}(s | s_0, a_0) = \mathbb{E}_{\Delta t \sim \text{Geom}(1-\gamma)} [\mathbb{P}[S_{t+\Delta t} | s_0, a_0; \Delta t; \pi]] \quad (4)$$

This decomposition of the value function has already been used in previous works based on the successor representation ([Dayan, 1993](#); [Barreto et al., 2016](#)) and, more recently,  $\gamma$ -models ([Janner et al., 2020](#)). We will use it to efficiently learn an implicit density ratio proportional to the state occupancy measure using contrastive learning.

<sup>1</sup> $\Delta(\mathcal{X})$  denotes the entire set of distributions over the space  $\mathcal{X}$ .

**Noise-contrastive estimation** Noise-contrastive estimation (NCE, [Gutmann and Hyvärinen, 2010](#)) spans a broad class of learning algorithms, at the core of which is negative sampling ([Mikolov et al., 2013](#)), i.e., learning an implicit metric space from positive and negative examples. Given reference samples, samples from a positive distribution (i.e., high similarity with reference points) and samples from a negative distribution (i.e., low similarity with reference points), contrastive learning methods learn an embedding where positive examples are located closer to the reference points than negative examples. One of the most well-known and commonly used NCE objectives is InfoNCE ([Oord et al., 2018](#)), which solves

$$\max_{\phi, \psi \in \Phi} \mathbb{E}_{x, y, y'} \left[ \log \frac{e^{\phi(x)^\top \psi(y)}}{\sum_{y' \in y \cup y'} e^{\phi(x)^\top \psi(y')}} \right] \quad (5)$$

over some hypothesis class  $\Phi : \{\phi : \mathcal{X} \rightarrow \mathcal{Z}\}$  for input space  $\mathcal{X}$ , latent space  $\mathcal{Z}$ ,  $x \sim \mathbb{P}(\mathcal{X})$ ,  $y \sim \mathbb{P}_{\text{positives}}(\mathcal{X})$  and  $y' \sim \mathbb{P}_{\text{negatives}}(\mathcal{X})$ . Contrastive learning has been widely studied in the static unsupervised/ supervised learning settings ([Hjelm et al., 2018](#); [Chen et al., 2020](#); [He et al., 2020](#)), as well as in reinforcement learning ([Kim et al., 2018](#); [Mazouze et al., 2020](#)) for learning state representations with desirable properties such as alignment and uniformity ([Wang and Isola, 2020](#)).

Solving Equation (5) for  $(\phi^*, \psi^*)$  yields a critic  $f : \mathcal{X} \times \mathcal{Y} \rightarrow \mathbb{R}$  which decomposes as  $f^*(x, y) = \phi^*(x)^\top \psi^*(y)$  and, at optimality<sup>2</sup>, captures the log-ratio of  $\mathbb{P}_{\text{positives}}(\mathcal{X})$  and  $\mathbb{P}_{\text{negatives}}(\mathcal{X})$ :

$$f^*(x, y) \propto \log \frac{\mathbb{P}[y|x]}{\mathbb{P}[y]}. \quad (6)$$

**Implicit dynamics models via NCE.** Various prior works ([Du and Mordatch, 2019](#); [Mazouze et al., 2020](#); [Nachum and Yang, 2021](#)) have studied the use of NCE to approximate a single-step dynamics model, where triplets  $(s_t, a_t, s_{t+1})$  have higher similarity than  $(s_t, a_t, s_{t' \neq t+1})$ , effectively defining positive and negative distributions over trajectory data. More recently, contrastive goal-conditioned RL ([Eysenbach et al., 2022](#)) used InfoNCE to condition the critic on goal states sampled from the replay buffer. These methods use asymmetric encoders, using  $\phi(s_t, a_t)$  and  $\psi(s_{t+\Delta t})$ , where positive samples of  $s_{t+\Delta t}$  are sampled from the discounted state occupancy measure for  $t \geq 0$ .

The conditional probability distribution of future states given the current state-action pair can be efficiently estimated using an implicit model trained via contrastive learning over positive and negative feature distributions, as shown in Equation (7).

$$\ell_{\text{InfoNCE}}(\phi, \psi) = \mathbb{E}_{s_t, a_t, \Delta t, \Delta t'} \left[ -\log \frac{e^{\phi(s_t, a_t)^\top \psi(s_{t+\Delta t})}}{\sum_{\Delta t' \in \Delta t \cup \Delta t'} e^{\phi(s_t, a_t)^\top \psi(s_{t+\Delta t'})}} \right] \quad (7)$$

Minimizing  $\ell_{\text{InfoNCE}}$  over trajectory data yields a critic which, at optimality, approximates the future discounted state occupancy measure up to a multiplicative term as per Equation (6),

$$f^*(s_t, a_t, s_{t+\Delta t}) \propto \log \frac{\mathbb{P}[s_{t+\Delta t} | s_t, a_t; \pi]}{\mathbb{P}[s_{t+\Delta t}; \pi]}. \quad (8)$$

Intuitively,  $f^*$  approximates a  $H$ -step dynamics model which has an implicit dependence on policy  $\pi$  that collected the training data. Ordinarily, training state-space models is hard when the dimensions are large, e.g. image-based domains. However, by using contrastive learning, we can learn this model without having to require it predict high-dimensional observations, as similarity is evaluated in a lower-dimensional latent space (observe that in Equation (7) the inner product is computed in  $\mathcal{Z}$ , whose dimension we control, instead of  $\mathcal{X}$ , which is specified externally). An apparent limitation of the approach is that the probability of future states  $s_{t+\Delta t}$  is recovered only up to a constant. However, it turns out that we can still use this model to get accurate estimates of the Q-values, as is described in the next section.

## 4 ESTIMATING AND MAXIMIZING RETURNS VIA CONTRASTIVE LEARNING

In this section, we show how NCE can be used to learn a quantity proportional to a value function, and how the later can be used in a policy iteration scheme.

<sup>2</sup>See [Ma and Collins \(2018\)](#) for exact derivation.

#### 4.1 ESTIMATING Q-VALUES USING THE CONTRASTIVE MODEL

As shown in Equation (3), the Q-function at  $(s_t, a_t)$  can be thought of as evaluating the reward function at states sampled from the discounted occupancy measure  $\mathbb{P}_{t:H}^\pi(s_t, a_t)$ . That is, to estimate a quantity akin to  $Q^\pi$ , we can first estimate the occupancy measure and take a weighted average of rewards over future states using the probabilities from the log-density ratio learned by the contrastive model. Precisely, Equation (3) corresponds to using an importance-weighted estimator, where an optimal critic that minimizes Equation (7) approximates the density ratio from Equation (8).

The critic itself can be trained using the occupancy measure formulation specified in Equation (4) over all state-action pairs in a given episode. However, Equation (4) needs to be re-adjusted to account for finite-horizon truncation of the geometric mass function presented in Definition 1.

**Definition 1 (Truncated distribution)** Let  $X$  be a random variable with distribution function  $F_X$ .  $Y$  is called the **truncated distribution** of  $X$  with support  $[m, M]$  s.t.  $0 < m < M$  if

$$\mathbb{P}[Y = y] = \frac{F_X(y - m) - F_X(y - 1 - m)}{F_X(M) - F_X(m)}, \quad y = m, m + 1, m + 2, \dots, M \quad (9)$$

We denote the special case of the truncated geometric distribution as  $\text{TruncGeom}(p, m, M)$ .

The contrastive objective to train the critic to approximate the discounted occupancy measure over a dataset  $\mathcal{D}$  is then

$$\ell_{\text{OM-InfoNCE}}(\phi, \psi) = \mathbb{E}_{\substack{s_t, a_t \sim \mathcal{D}, \\ \Delta t \sim \text{TruncGeom}(1-\gamma, t, H), \\ \Delta t' \sim \text{TruncGeom}(1-\gamma, t' \neq t, H)}} \left[ -\log \frac{e^{\phi(s_t, a_t)^\top \psi(s_{t+\Delta t})}}{\sum_{\Delta t' \in \Delta t \cup \Delta t'} e^{\phi(s_t, a_t)^\top \psi(s_{t+\Delta t'})}} \right] \quad (10)$$

It is possible that multiple optimal critics exist s.t. the multiplicative proportionality constant depends on the action. To avoid this, we adopt a similar approach as Eysenbach et al. (2022) and introduce a regularization term over the partition function, making the critic training objective be

$$\ell_{\text{Critic}} = \ell_{\text{OM-InfoNCE}} + \lambda_{\text{Partition}} \mathbb{E}_{s_t, a_t, \Delta t, \Delta t'} [(\log \sum_{\Delta t' \in \Delta t \cup \Delta t'} e^{\phi(s_t, a_t)^\top \psi(s_{t+\Delta t'})})^2] \quad (11)$$

Now, suppose we found an optimal critic  $f$ . Combining Equation (4) with Definition 1, we obtain the following form of the Q-function for an optimal critic  $f$  which minimizes Equation (7):

$$\begin{aligned} Q_{\text{NCE}}(s_t, a_t) &= \sum_{\Delta t=1}^{\infty} \gamma^{\Delta t-1} \int_{s_{t+\Delta t}} r(s_{t+\Delta t}) \mathbb{P}[s_{t+\Delta t} | s_t, a_t; \pi] ds_{t+\Delta t} \\ &\propto \frac{1}{1-\gamma} \mathbb{E}_{\Delta t \sim \text{TruncGeom}(1-\gamma, t, H)} \left[ \int_{s_{t+\Delta t}} r(s_{t+\Delta t}) e^{f(s_t, a_t, s_{t+\Delta t})} \mathbb{P}[s_{t+\Delta t}; \pi] ds_{t+\Delta t} \right] \\ &= \frac{1}{1-\gamma} \mathbb{E}_{\Delta t \sim \text{TruncGeom}(1-\gamma, t, H)} [\mathbb{E}_{\mathbb{P}_{t+\Delta t}^\pi} [r(s_{t+\Delta t}) e^{f(s_t, a_t, s_{t+\Delta t})}]] \end{aligned} \quad (12)$$

Here, the offset  $\Delta t$  is a random variable sampled from  $\text{TruncGeom}(1-\gamma, t, H)$  where  $H$  is the horizon of the MDP<sup>3</sup>. We can also show that  $Q(s, a) < Q(s, a') \implies Q_{\text{NCE}}(s, a) < Q_{\text{NCE}}(s, a')$  for all  $s \in \mathcal{S}$  and  $a, a' \in \mathcal{A}$ , which makes the contrastive Q-values suitable for policy evaluation. However, we do not, in general, expect  $Q_{\text{NCE}}$  to recover the optimal Q function, as the recovered Q-values are on-policy with respect to  $\pi$ .

#### 4.2 EFFICIENT ESTIMATION USING RANDOM FOURIER FEATURES

A major issue with using  $Q_{\text{NCE}}$  out-of-the-box is that it is computationally expensive, requiring evaluation of the inner product  $\phi(s_t, a_t)^\top \psi(s_{t+\Delta t})$  with a large number of future states and hence multiple forward passes through  $\psi$ . The underlying cause of this computational overhead is the RBF

<sup>3</sup>While using the truncated geometric distribution makes Equation (12) proportional to the true value function, the relation becomes an equality in the infinite-horizon case since  $\lim_{H \rightarrow \infty} 1 - (1-p)^{H-m} = 1$ .

kernel term  $e^{\phi(s_t, a_t)^\top \psi(s_{t+\Delta t})}$ . If we instead used a linear kernel, the constant term  $\phi(s_t, a_t)$  would be factored out, and we could separately keep track of reward-weighted future expected features. This would (1) reduce the computational complexity of  $N$  actor updates over  $\mathcal{D}$  from  $\mathcal{O}(|\mathcal{D}|N)$  to  $\mathcal{O}(|\mathcal{D}| + N)$  and (2) reduce the variance of the representation if averaging features of future states using exponential moving average. It turns out that the RBF kernel can be approximately linearized by using random Fourier features (Rahimi and Recht, 2007; Nachum and Yang, 2021).

**Lemma 1** *Let  $\mathbf{x}, \mathbf{y} \in \mathbb{R}^d$  be unit vectors, and let  $F_{\mathbf{W}, \mathbf{b}}(\mathbf{x}) = \sqrt{\frac{2e}{d}} \cos(\mathbf{W}\mathbf{x} + \mathbf{b})$  where  $\mathbf{W} \sim \text{Normal}(0, \mathbf{I})$  and  $\mathbf{b} \sim \text{Uniform}(0, 2\pi)$ . Then,  $\mathbb{E}[F_{\mathbf{W}, \mathbf{b}}(\mathbf{x})^\top F_{\mathbf{W}, \mathbf{b}}(\mathbf{y})] = e^{\mathbf{x}^\top \mathbf{y}}$ .*

Lemma 1 is a straightforward modification of the result from Rahimi and Recht (2007) and allows us to reduce the RBF kernel to an expectation over  $d$ -dimensional random feature vectors:

$$\begin{aligned} Q_{\text{NCE}}(s_t, a_t) &= \frac{1}{1-\gamma} \mathbb{E}_{\Delta t \sim \text{TruncGeom}(1-\gamma, t, H)} [\mathbb{E}_{\mathbb{P}(s_{t+\Delta t}; \pi)} [e^{\phi(s_t, a_t)^\top \psi(s_{t+\Delta t})} r(s_{t+\Delta t})]] \\ &= \frac{1}{1-\gamma} F_{\mathbf{W}, \mathbf{b}}(\phi(s_t, a_t))^\top \mathbb{E}_{\Delta t \sim \text{TruncGeom}(1-\gamma, t, H)} [\mathbb{E}_{\mathbb{P}(s_{t+\Delta t}; \pi)} [F_{\mathbf{W}, \mathbf{b}}(\psi(s_{t+\Delta t})) r(s_{t+\Delta t})]] \\ &= \frac{1}{1-\gamma} F_{\mathbf{W}, \mathbf{b}}(\phi(s_t, a_t)) \xi(\pi) \end{aligned} \tag{13}$$

The advantage of using the RFF approximation is that it allows us to split the exponential term inside the expectation and separately keep track of the policy-dependent, reward-weighted future state probability term, while the state-action dependence term is learned online. Specifically, we keep track of  $\xi(\pi)$  via an exponential-moving average during the entire duration of training<sup>4</sup>.

### 4.3 LEARNING THE POLICY

Once the policy evaluation phase completes and we have an estimate  $Q_{\text{NCE}}$ , we optimize a policy to maximize the returns predicted by this Q-value. We can decode the policy by minimizing its Kullback-Leibler divergence to the Boltzmann Q-value distribution (see Haarnoja et al. (2018)), which can be efficiently done by minimizing the following objective:

$$\ell_{\text{Policy}}(\theta) = \mathbb{E}_{s_t \sim \mathcal{D}} \left[ \text{D}_{\text{KL}} \left( \pi_\theta(s_t) \left\| \frac{e^{Q(s_t, \cdot)/\tau}}{\int_{a \in \mathcal{A}} e^{Q(s_t, a)/\tau} da} \right. \right) \right]. \tag{14}$$

Note that in discrete action spaces, minimizing Equation (14) leads to a soft version of the greedy policy decoding  $\pi_{\text{greedy}}(s) = \arg \max_{a \in \mathcal{A}} Q_{\text{NCE}}(s, a)$  for  $s \in \mathcal{S}$ . In practice, we approximate the KL term in Equation (14) using  $N_a$  Monte-Carlo action samples.

Decoding  $\pi$  in such a way can lead to sampling out-of-distribution actions in regions where the Q-function might be inaccurate due to poor dataset coverage. To mitigate this issue, we follow prior work (Cobbe et al., 2021; Zhao et al., 2021; Schwarzer et al., 2021) and add a behavior cloning term which prevents the new policy from straying too far away from the data:

$$\ell_{\text{BC}}(\theta) = \mathbb{E}_{a, s \sim \mathcal{D}} [\log \pi_\theta(a|s)] + \tau \mathbb{E}_{s \sim \mathcal{D}} [\mathcal{H}(\pi_\theta(s))]. \tag{15}$$

for entropy estimator  $\mathcal{H}(\pi(s)) = -\mathbb{E}_{a \sim \pi(s)} [\log \pi(a|s)]$ . We add this extra loss to  $\ell_{\text{Policy}}$  to learn a policy  $\pi$  which prioritizes high Q-values that are well-supported by the offline dataset  $\mathcal{D}$ . Thus, the final policy optimization objective becomes

$$\ell_{\text{Policy}}(\theta) = \ell_{\text{Policy}}(\theta) + \lambda_{\text{BC}} \ell_{\text{BC}}(\theta). \tag{16}$$

The policy found by minimizing  $\ell_{\text{Policy}}$  has, on average, non-decreasing returns, as per Lemma 2.

**Lemma 2 (Contrastive policy improvement)** *Let  $\mu$  be a policy and let  $Q_{\text{NCE}}^\mu = \min_{\phi, \psi \in \Phi} \mathbb{E}_{\mathcal{D}^\mu} [\ell_{\text{Critic}}(\phi, \psi)]$ . If*

$$\pi(s) = \arg \min_{\pi \in \Pi} D_{\text{KL}} \left( \pi(s) \left\| \frac{e^{Q_{\text{NCE}}^\mu(s_t, \cdot)/\tau}}{\int_{a \in \mathcal{A}} e^{Q_{\text{NCE}}^\mu(s_t, a)/\tau} da} \right. \right) \tag{17}$$

*then  $Q^\pi(s, a) \geq Q^\mu(s, a)$  for all  $(s, a) \in \mathcal{D}^\mu$ .*

<sup>4</sup>This idea can be adapted to online learning settings as well by clipping policy improvement steps so that  $\xi$  doesn't change too fast under newly collected data.

**Algorithm 1:** Contrastive Value Learning (CVL)

---

**Input** : Dataset  $\mathcal{D} \sim \mu, \psi, \phi$  networks, temperature parameter  $\tau$ , exponential moving average parameter  $\beta$

- 1 **for** epoch  $j = 1, 2, \dots, J$  **do**
- 2     **for** minibatch  $\mathcal{B} \sim \mathcal{D}$  **do**
- 3         /\* Update density ratio estimator using Equation (11) \*/
- 3         Update  $\phi^{(j+1)}, \psi^{(j+1)}$  using  $\nabla_{\phi, \psi} \ell_{\text{Critic}}(\phi^{(j)}, \psi^{(j)})$  ;
- 4         /\* Estimate the contrastive Q-function \*/
- 4          $Q(s, a) \leftarrow$  Equation (13) if using RFF, otherwise Equation (12);
- 4         /\* Decode policy from Q-function using Equation (16) \*/
- 5         Update  $\pi_{\theta}$  using  $\nabla_{\theta} \{\ell_{\text{Policy}}(\theta)\}$  ;
- 5         /\* Update future state encoder using EMA \*/
- 6          $\psi^{(j+1)} \leftarrow \beta \psi^{(j+1)} + (1 - \beta) \psi^{(j)}$  ;
- 6         /\* Update future state state features weighted by rewards \*/
- 7          $\xi^{(j+1)} \leftarrow \psi^{(j+1)} \cdot \mathcal{B}[r_{t+\Delta t}]$  ;

---

The proof of Lemma 2 is located in Section 6.2. Specifically, Lemma 2 tells us that using CVL as a surrogate Q-function corresponds to one step of conservative policy improvement, where  $\pi$  satisfies soft constraints of Equation (14) and small  $\mathbb{E}_{\mathcal{D}^{\mu}} [D_{\text{KL}}(\pi(s) || \mu(s))]$  via the BC term.

#### 4.4 PRACTICAL IMPLEMENTATION

We now present our complete method, which can be viewed as an actor-critic method for offline RL. We learn the critic via contrastive learning (Equation (11)) and learn the policy via Equation (16). We will interleave these steps in most of our experiments, but experiments in Section 5.1 show that the critic can be pretrained e.g. in the presence of unlabeled data from related tasks. We summarize the method in Algorithm 1.

#### 4.5 INTERPRETATIONS AND CONNECTIONS WITH PRIOR WORK

The main distinction between Contrastive Value Learning and prior works consists specifically in representing the Q-values in a two-step decomposition: the Q-value is represented as an occupancy measure weighted by the reward signal; the occupancy measure itself is represented using a powerful likelihood-based model parameterized using an implicit function. Decoupling the learning of the occupancy measure from reward maximization allows, among others, for efficient pretraining strategies on unlabeled data, i.e. trajectory data without reward information, and can be used to learn provably optimal state representations for *any* reward function (Touati and Ollivier, 2021). While CVL is similar in spirit to the successor representation (Dayan, 1993; Barreto et al., 2016), the occupancy measure learned by CVL is much richer than that of SR, as it captures the entire distribution over future states instead of only the first moment. Another method,  $\gamma$ -models (Janner et al., 2020), is closely related to CVL, but uses a surrogate single-step TD objective to learn the occupancy measure, similarly to C-learning (Eysenbach et al., 2020).

## 5 EXPERIMENTS

Our experiments aim to answer three questions. First, we study how CVL compares with baseline approaches on a large benchmark of state-based tasks. Our second set of experiments look at image-based tasks, testing the hypothesis that CVL scales to these tasks more effectively than the baselines. We conclude with ablation experiments. Our main point of comparison will be a high-performing offline RL method, CQL (Kumar et al., 2020). While CVL learns an implicit model, that model is structurally much more similar to value-based RL methods than model-based methods, motivating our comparison to a value-based baseline (CQL). We will also include behavioral cloning as a baseline.

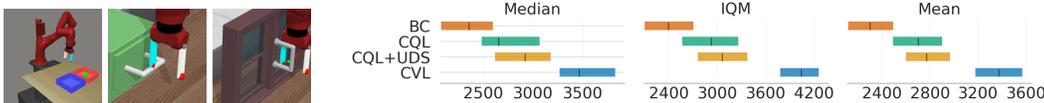


Figure 2: **Metaworld benchmark.** (Left) We evaluate CVL on 50 tasks from Metaworld, a subset of which are shown here. (Right) Compared with three offline RL baselines, CVL achieves statistically-significant improvements in offline performance. Results are reported over 5 random seeds.

**Metaworld.** We first test our approach on the complex MetaWorld benchmark (Yu et al., 2020a), which consists of 50 robotic manipulation tasks such as open a door, pick up an object, reach a certain area of the table, executed by a robotic arm (see Figure 2 (left)). This domain is an ideal testbed for CVL, as it allows for both full-state and image-based experiments, has a dense and informative reward function thus decoupling the problem of representation learning from exploration, and is challenging for model-free methods which leaves room for improvement. While the original MetaWorld domain has been used to evaluate online RL agents, we create an *ad hoc* dataset suitable for offline learning. To do so, we train Soft Actor-Critic (Haarnoja et al., 2018) from full states on each of the 50 tasks separately for 500k frames, and save the resulting replay buffer, which forms the training dataset. As shown in Figure 2 (right), CVL manages to considerably improve upon strong baselines such as behavior cloning, CQL and CQL with UDS (Yu et al., 2022)<sup>5</sup>. We report the results on all tasks of the MetaWorld suite over 5 random seeds, according to the aggregation methodology proposed by Agarwal et al. (2021b). Per-environment scores are available in Table 4.

**Image-based experiments** Our working hypothesis is that contrastive formulation of the value function acts in itself as a pre-training mechanism via the prism of representation learning. For this reason, we conduct further experiments on 4 image-based tasks from the MetaWorld suite (similarly to full-states, the dataset was obtained from the SAC replay buffer trained on rendered images).

Table 1: **Offline RL with Images.** We compare CVL to baselines on four offline, image-based tasks from MetaWorld offline image-based tasks. Average  $\pm$  std. dev. are shown for 5 random seeds.

Task	BC	CQL	CVL
door-close	571 $\pm$ 9.9	4249 $\pm$ 269.9	<b>4480 <math>\pm</math> 305.1</b>
door-open	178 $\pm$ 4.0	2099 $\pm$ 0.9	<b>3389 <math>\pm</math> 76.6</b>
drawer-close	2414 $\pm$ 1736.5	<b>3964 <math>\pm</math> 1634.9</b>	2177 $\pm$ 1679.5
drawer-open	1030 $\pm$ 104.2	820 $\pm$ 56.0	<b>2543 <math>\pm</math> 115.0</b>

Results presented in Table 1 show that CVL is also able to learn meaningful Q-values and achieve good empirical performance on hard image-based tasks.

### 5.1 ABLATION EXPERIMENTS.

**When is pretraining the model useful?** In theory, the model can be pretrained on the data from other tasks, however, we do not always expect this to help (e.g., when the pretraining tasks are very different). We ran an experiment to test this capability. The results, shown in Fig. 3, show that pretraining sometimes speeds up learning. In particular, we observe that pretraining is effective when the pretraining tasks are similar to the target task and contain a diverse set of state-action pairs.

**How reliable is the  $Q_{NCE}$  approximation?** Given that contrastive Q-values are proportional to the true Q-function, a natural question to ask is how good is  $Q_{NCE}$  at capturing the topology of  $Q$ ? First, we conduct an ablation demonstrating how linearizing the RBF kernel via random Fourier features provides a performance gain on the offline MetaWorld tasks Figure 4. Specifically, we hypothesize that this is due to the reduced variance of the RFF Q-value estimator which keeps track of future reward-weighted state features using a rolling average.

Next, we qualitatively assess the similarity between contrastive and true Q-values on the continuous Mountain Car environment (Moore, 1990) by first pre-training SAC online on the task and then

<sup>5</sup>For CQL+UDS, we combine all data from the current task with unlabeled data from related tasks with rewards set to 0. In the absence of related tasks, we pre-train the critic on the current task with 0 rewards.

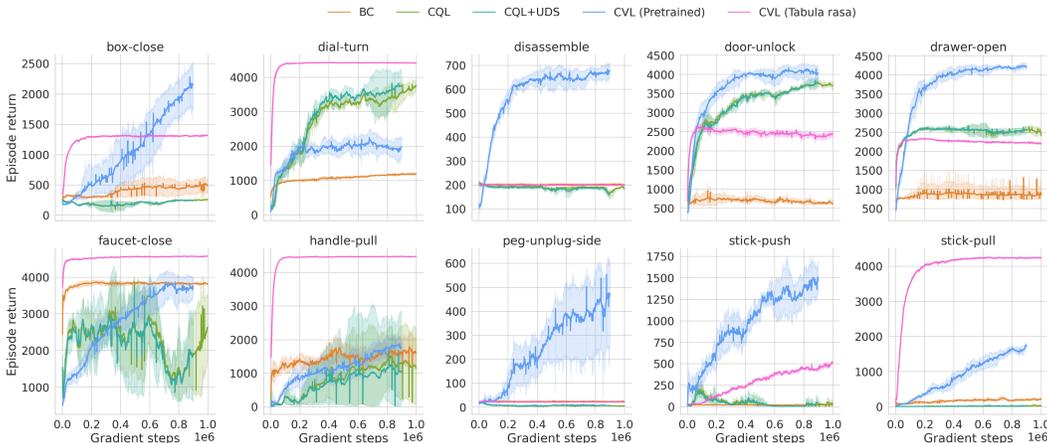


Figure 3: **Offline Learning Curves for Metaworld.** Episode return curves as a function of gradient steps taken during training on 10 random MetaWorld tasks; curves show mean  $\pm$  standard deviation. Pretraining the reward-free occupancy measure on related tasks allows CVL to outperform baseline approaches and even CVL trained *tabula rasa*.

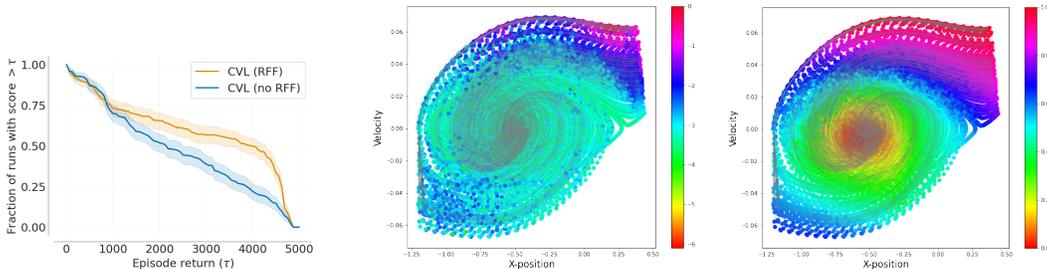


Figure 4: CVL with RFF (orange) performs slightly better than without RFF (blue).

Figure 5: **Visualizing the estimated Q-values.** (Left) Normalized  $\log Q_{NCE}$  learned by CVL offline on the Mountain Car environment. (Right) Normalized  $Q$  learned by online SAC on the same environment.

fitting CVL to the data from SAC’s replay buffer. Figure 5 (left) shows the contrastive Q-values on a log-scale, evaluated on trajectories from the SAC replay; for comparison, we also show the Q-values learned by online SAC in Figure 5 (right). Note that the value function learned by CVL conserves the same topology as the true value function, up to a multiplicative rescaling.

## 6 DISCUSSION

This paper presented an RL algorithm that learns a contrastive model of the world, and uses that model to obtain Q values by estimating the likelihood of visiting future states. Our experiments demonstrate that this approach can effectively solve a large number of offline RL tasks, including from image-based observations. Our pretraining results hinted that CVL can be pretrained on datasets from other tasks, and we are excited to pretrain our model on datasets of increasing size.

**Limitations.** One limitation of our approach is that it corresponds to a single step of policy improvement. This limitation might be lifted by training the contrastive model using a temporal difference update for the contrastive model (Eysenbach et al., 2020; Blier et al., 2021). A second limitation is that the RFF approximation can be poor when the feature dimension is small. We tried to train the contrastive model using non-exponentiated features (akin to HaoChen et al. (2021)), but failed to achieve satisfactory results. Figuring out how to effectively train these spectral (i.e., non-exponentiated) models remains an important question.

## REPRODUCIBILITY STATEMENT

We ensure reproducibility of our method via a) releasing the offline MetaWorld dataset that we used for our main results to allow the community to conduct further research on this domain upon publication, b) releasing the code used to obtain our results upon publication and c) detailing the hyperparameters and design choices for the implementation of CVL and the computational resources used for our experiments in Section 6.1.

## REFERENCES

- R. Agarwal, M. C. Machado, P. S. Castro, and M. G. Bellemare. Contrastive behavioral similarity embeddings for generalization in reinforcement learning. *arXiv preprint arXiv:2101.05265*, 2021a.
- R. Agarwal, M. Schwarzer, P. S. Castro, A. C. Courville, and M. Bellemare. Deep reinforcement learning at the edge of the statistical precipice. *Advances in neural information processing systems*, 34:29304–29320, 2021b.
- A. Argenson and G. Dulac-Arnold. Model-based offline planning. *arXiv preprint arXiv:2008.05556*, 2020.
- J. L. Ba, J. R. Kiros, and G. E. Hinton. Layer normalization. *arXiv preprint arXiv:1607.06450*, 2016.
- A. Barreto, W. Dabney, R. Munos, J. J. Hunt, T. Schaul, H. Van Hasselt, and D. Silver. Successor features for transfer in reinforcement learning. *arXiv preprint arXiv:1606.05312*, 2016.
- L. Blier, C. Tallec, and Y. Ollivier. Learning successor states and goal-dependent values: A mathematical viewpoint. *arXiv preprint arXiv:2101.07123*, 2021.
- T. Chen, S. Kornblith, M. Norouzi, and G. Hinton. A simple framework for contrastive learning of visual representations. In *International conference on machine learning*, pages 1597–1607. PMLR, 2020.
- K. W. Cobbe, J. Hilton, O. Klimov, and J. Schulman. Phasic policy gradient. In *International Conference on Machine Learning*, pages 2020–2027. PMLR, 2021.
- P. Dayan. Improving generalization for temporal difference learning: The successor representation. *Neural Computation*, 5(4):613–624, 1993.
- Y. Du and I. Mordatch. Implicit generation and modeling with energy based models. *Advances in Neural Information Processing Systems*, 32, 2019.
- L. Espeholt, H. Soyer, R. Munos, K. Simonyan, V. Mnih, T. Ward, Y. Doron, V. Firoiu, T. Harley, I. Dunning, et al. Impala: Scalable distributed deep-rl with importance weighted actor-learner architectures. In *International Conference on Machine Learning*, pages 1407–1416. PMLR, 2018.
- B. Eysenbach, R. Salakhutdinov, and S. Levine. C-learning: Learning to achieve goals via recursive classification. *arXiv preprint arXiv:2011.08909*, 2020.
- B. Eysenbach, T. Zhang, R. Salakhutdinov, and S. Levine. Contrastive learning as goal-conditioned reinforcement learning. *arXiv preprint arXiv:2206.07568*, 2022.
- S. Fujimoto, D. Meger, and D. Precup. Off-policy deep reinforcement learning without exploration. In *International Conference on Machine Learning*, pages 2052–2062. PMLR, 2019.
- M. Gutmann and A. Hyvärinen. Noise-contrastive estimation: A new estimation principle for unnormalized statistical models. In *Proceedings of the thirteenth international conference on artificial intelligence and statistics*, pages 297–304. JMLR Workshop and Conference Proceedings, 2010.
- T. Haarnoja, A. Zhou, P. Abbeel, and S. Levine. Soft actor-critic: Off-policy maximum entropy deep reinforcement learning with a stochastic actor. In *International conference on machine learning*, pages 1861–1870. PMLR, 2018.

- J. Z. HaoChen, C. Wei, A. Gaidon, and T. Ma. Provable guarantees for self-supervised deep learning with spectral contrastive loss. *Advances in Neural Information Processing Systems*, 34:5000–5011, 2021.
- K. He, H. Fan, Y. Wu, S. Xie, and R. Girshick. Momentum contrast for unsupervised visual representation learning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 9729–9738, 2020.
- R. D. Hjelm, A. Fedorov, S. Lavoie-Marchildon, K. Grewal, P. Bachman, A. Trischler, and Y. Bengio. Learning deep representations by mutual information estimation and maximization. *arXiv preprint arXiv:1808.06670*, 2018.
- G. Huang, Z. Liu, L. Van Der Maaten, and K. Q. Weinberger. Densely connected convolutional networks. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 4700–4708, 2017.
- M. Janner, I. Mordatch, and S. Levine. Generative temporal difference learning for infinite-horizon prediction. *arXiv preprint arXiv:2010.14496*, 2020.
- R. Kidambi, A. Rajeswaran, P. Netrapalli, and T. Joachims. Morel: Model-based offline reinforcement learning. *arXiv preprint arXiv:2005.05951*, 2020.
- H. Kim, J. Kim, Y. Jeong, S. Levine, and H. O. Song. Emi: Exploration with mutual information. *arXiv preprint arXiv:1810.01176*, 2018.
- A. Kumar, A. Zhou, G. Tucker, and S. Levine. Conservative q-learning for offline reinforcement learning. *arXiv preprint arXiv:2006.04779*, 2020.
- Z. Ma and M. Collins. Noise contrastive estimation and negative sampling for conditional models: Consistency and statistical efficiency. *arXiv preprint arXiv:1809.01812*, 2018.
- B. Mazouze, R. T. d. Combes, T. Doan, P. Bachman, and R. D. Hjelm. Deep reinforcement and infomax learning. *Neural Information Processing Systems*, 2020.
- T. Mikolov, K. Chen, G. Corrado, and J. Dean. Efficient estimation of word representations in vector space. *arXiv preprint arXiv:1301.3781*, 2013.
- A. W. Moore. Efficient memory-based learning for robot control. 1990.
- R. Munos. Error bounds for approximate policy iteration. In *ICML*, volume 3, pages 560–567, 2003.
- O. Nachum and M. Yang. Provable representation learning for imitation with contrastive fourier features. *Advances in Neural Information Processing Systems*, 34:30100–30112, 2021.
- A. v. d. Oord, Y. Li, and O. Vinyals. Representation learning with contrastive predictive coding. *arXiv preprint arXiv:1807.03748*, 2018.
- A. Rahimi and B. Recht. Random features for large-scale kernel machines. *Advances in neural information processing systems*, 20, 2007.
- M. Schwarzer, N. Rajkumar, M. Noukhovitch, A. Anand, L. Charlin, D. Hjelm, P. Bachman, and A. Courville. Pretraining representations for data-efficient reinforcement learning. *arXiv preprint arXiv:2106.04799*, 2021.
- A. Srinivas, M. Laskin, and P. Abbeel. Curl: Contrastive unsupervised representations for reinforcement learning. *International Conference on Machine Learning*, 2020.
- A. Touati and Y. Ollivier. Learning one representation to optimize all rewards. *Advances in Neural Information Processing Systems*, 34, 2021.
- T. Wang and P. Isola. Understanding contrastive representation learning through alignment and uniformity on the hypersphere. In *International Conference on Machine Learning*, pages 9929–9939. PMLR, 2020.

- Y. Wu, G. Tucker, and O. Nachum. Behavior regularized offline reinforcement learning. *arXiv preprint arXiv:1911.11361*, 2019.
- T. Yu, D. Quillen, Z. He, R. Julian, K. Hausman, C. Finn, and S. Levine. Meta-world: A benchmark and evaluation for multi-task and meta reinforcement learning. In *Conference on Robot Learning*, pages 1094–1100. PMLR, 2020a.
- T. Yu, G. Thomas, L. Yu, S. Ermon, J. Y. Zou, S. Levine, C. Finn, and T. Ma. Mopo: Model-based of-line policy optimization. *Advances in Neural Information Processing Systems*, 33:14129–14142, 2020b.
- T. Yu, A. Kumar, R. Rafailov, A. Rajeswaran, S. Levine, and C. Finn. Combo: Conservative offline model-based policy optimization. *Advances in neural information processing systems*, 34:28954–28967, 2021.
- T. Yu, A. Kumar, Y. Chebotar, K. Hausman, C. Finn, and S. Levine. How to leverage unlabeled data in offline reinforcement learning. *arXiv preprint arXiv:2202.01741*, 2022.
- Y. Zhao, R. Boney, A. Ilin, J. Kannala, and J. Pajarinen. Adaptive behavior cloning regularization for stable offline-to-online reinforcement learning. 2021.

## APPENDIX

### REPRODUCIBILITY CHECKLIST

1. For all authors...
  - (a) Do the main claims made in the abstract and introduction accurately reflect the paper’s contributions and scope? **Yes**
  - (b) Did you describe the limitations of your work? **See end of Section 4**
  - (c) Did you discuss any potential negative societal impacts of your work? **N/A**
  - (d) Have you read the ethics review guidelines and ensured that your paper conforms to them? **Yes**
2. If you are including theoretical results...
  - (a) Did you state the full set of assumptions of all theoretical results? **Yes, they are outlined in the lemmas.**
  - (b) Did you include complete proofs of all theoretical results? **See Appendix Section 6.2**
3. If you ran experiments...
  - (a) Did you include the code, data, and instructions needed to reproduce the main experimental results (either in the supplemental material or as a URL)? **Code will be release upon publication, along with the offline Metaworld dataset.**
  - (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they were chosen)? **See Appendix Section 6.1**
  - (c) Did you report error bars (e.g., with respect to the random seed after running experiments multiple times)? **Yes, we include multiple various uncertainty estimates and measures of central tendency, e.g as outlined in Agarwal et al. (2021b) over 5 random replicates (i.e. seeds)**
  - (d) Did you include the total amount of compute and the type of resources used (e.g., type of GPUs, internal cluster, or cloud provider)? **See Appendix Section 6.1**
4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets...
  - (a) If your work uses existing assets, did you cite the creators? **N/A**
  - (b) Did you mention the license of the assets? **N/A**
  - (c) Did you include any new assets either in the supplemental material or as a URL? **Due to the large size of the offline Metaworld dataset, we will release the code required to re-generate it exactly upon publication.**
  - (d) Did you discuss whether and how consent was obtained from people whose data you’re using/curating? **N/A**

- (e) Did you discuss whether the data you are using/curating contains personally identifiable information or offensive content? **N/A**
5. If you used crowdsourcing or conducted research with human subjects...
- (a) Did you include the full text of instructions given to participants and screenshots, if applicable? **N/A**
- (b) Did you describe any potential participant risks, with links to Institutional Review Board (IRB) approvals, if applicable? **N/A**
- (c) Did you include the estimated hourly wage paid to participants and the total amount spent on participant compensation? **N/A**

## 6.1 EXPERIMENTAL DETAILS

**Model architecture** All algorithms (baselines as well as CVL) were based on a common architecture, where an encoder (IMPALA (Espeholt et al., 2018) for image data and two layer DenseNet MLP (Huang et al., 2017) for full-states) generated state features which, combined with actions gave rise to the Q-value and the policy (we used a diagonal Gaussian policy with a Tanh bijector, as is common for continuous control tasks). The main difference of CVL with the baselines is that the critic is defined implicitly via the dot-product of current state-action features passed through one encoder, and future state features passed into a separate DenseNet. The output of both encoders was optionally normalized using  $\ell_2$  norm. All methods had a LayerNorm layer (Ba et al., 2016) in between each linear layer to ensure proper feature scaling.

Hyperparameter	Value
Learning rate	$3 \times 10^{-4}$
Batch size (all but CVL)	512
Discount factor	0.99
Framestack	No
Max gradient norm	100
MLP structure	$256 \times 256$ DenseNet
Encoder (full-state)	$256 \times 256$ DenseNet MLP
Encoder (pixels)	IMPALA
Add LayerNorm in between all layers	Yes

Table 2: Hyperparameters that are consistent between methods.

Hyperparameter	Value
CVL	
Batch size	$H$
Number of future action samples $N_a$	10
InfoNCE temperature	1
Partition function coefficient $\lambda_{\text{partition}}$	0.001
BC coefficient $\lambda_{\text{BC}}$	0 (Mountain Car), 0.1 (rest)
RFF	Yes
$\ell_2$ -normalize MLP outputs	Yes
CQL	
Regularization coefficient	1
BC	
Entropy regularization coefficient	0.1

Table 3: Hyperparameters that are different between methods.

All experiments were run on the equivalent of 8 V100 GPUs with 64 Gb of RAM and 8 CPUs.

**Dataset composition** The offline MetaWorld dataset was constructed by first pre-training SAC on all 50 tasks from full-states for 500k environment interactions. The replay buffer at the end of the training was then used as training dataset for BC, CQL, CQL+UDS and CVL. An identical approach was used to construct the image-based MetaWorld datasets and the Mountain Car dataset.

**Pretraining setup** When pretraining CVL, we first optimize the critic on unlabeled data from dataset for all the semantically related tasks, i.e. tasks which belong to the same domain, and then finetune both the critic and the policy on reward-labeled data from the target task. Semantically related tasks in MetaWorld are easily identifiable by their domain name, e.g. `drawer-open` and `drawer-close` belong to the `drawer` domain. We use a similar approach when pretraining CQL+UDS, where we perform TD updates with all rewards equal to 0 during the pretraining phase.

## 6.2 PROOFS

### Proof 1 (Random Fourier features approximation, Lemma 1)

For unit vectors  $x, y \in \mathbb{R}^d$ ,  $d > 0$ ,

$$\begin{aligned} \left( \sqrt{\frac{2}{d}} \cos(Wx + b) \right)^\top \left( \sqrt{\frac{2}{d}} \cos(Wy + b) \right) &= \exp\{-\|x - y\|_2^2/2\} \\ &= \exp\{-(\|x\|_2^2 - 2x^\top y + \|y\|_2^2)/2\} \\ &= \exp(-(2 - 2x^\top y)/2) \\ &= e^{x^\top y - 1} \\ &= \frac{e^{x^\top y}}{e} \end{aligned} \quad (18)$$

by re-arranging the terms in the result from [Rahimi and Recht \(2007\)](#). Therefore,

$$e^{x^\top y} = \left( \sqrt{\frac{2e}{d}} \cos(Wx + b) \right)^\top \left( \sqrt{\frac{2e}{d}} \cos(Wy + b) \right) \quad (19)$$

**Proof 2 (CVL induces a single-step of policy improvement, Lemma 2)** *Since, for the optimal critic  $f^*$ ,*

$$e^{f^*(s_t, a_t, s_{t+\Delta t})} \propto \frac{\mathbb{P}[s_{t+\Delta t} | s_t, a_t; \mu]}{\mathbb{P}[s_{t+\Delta t}; \mu]}. \quad (20)$$

*point-wise for every  $(s_t, a_t, s_{t+\Delta t}) \in \mathcal{D}^\mu$ , then, for  $\alpha > 0$ ,*

$$e^{f^*(s_t, a_t, s_{t+\Delta t})} = \alpha \frac{\mathbb{P}[s_{t+\Delta t} | s_t, a_t; \mu]}{\mathbb{P}[s_{t+\Delta t}; \mu]}. \quad (21)$$

*Now, the following relation holds using the previous result*

$$\begin{aligned} Q_{NCE}^\mu(s_t, a_t) &= \frac{1}{1 - \gamma} \mathbb{E}_{\Delta t \sim \text{TruncGeom}(1 - \gamma, t, H)} [\mathbb{E}_{\mathbb{P}_{t+\Delta t}^\mu} [r(s_{t+\Delta t}) e^{f^*(s_t, a_t, s_{t+\Delta t})}]] \\ &= \frac{\alpha}{1 - \gamma} \mathbb{E}_{\Delta t \sim \text{TruncGeom}(1 - \gamma, t, H)} \left[ \int_{s_{t+\Delta t}} r(s_{t+\Delta t}) \mathbb{P}[s_{t+\Delta t} | s_t, a_t; \mu] ds_{t+\Delta t} \right] \\ &= \alpha Q^\mu(s_t, a_t) \end{aligned} \quad (22)$$

*Using this relation yields*

$$\frac{e^{Q_{NCE}^\mu(s_t, a_t)/\tau}}{\int_{a \in \mathcal{A}} e^{Q_{NCE}^\mu(s_t, a)/\tau} da} = \frac{e^{\alpha Q^\mu(s_t, a_t)/\tau}}{\int_{a \in \mathcal{A}} e^{\alpha Q^\mu(s_t, a)/\tau} da} = \frac{e^{Q^\mu(s_t, a_t)/\tau}}{\int_{a \in \mathcal{A}} e^{Q^\mu(s_t, a)/\tau} da} \quad (23)$$

*It follows that*

$$\begin{aligned} \arg \min_{\pi \in \Pi} D_{KL} \left( \pi(s_t) \left\| \frac{e^{Q_{NCE}^\mu(s_t, \cdot)/\tau}}{\int_{a \in \mathcal{A}} e^{Q_{NCE}^\mu(s_t, a)/\tau} da} \right\| \right) &= \arg \min_{\pi \in \Pi} D_{KL} \left( \pi(s_t) \left\| \frac{e^{Q^\mu(s_t, \cdot)/\tau}}{\int_{a \in \mathcal{A}} e^{Q^\mu(s_t, a)/\tau} da} \right\| \right) \\ &= \pi(s_t) \end{aligned} \quad (24)$$

Now, we invoke Lemma 2 from [Haarnoja et al. \(2018\)](#) by using the equivalence of the policy decoded from contrastive  $Q$ -values to the policy found by soft policy iteration, which concludes the proof.

### 6.3 ADDITIONAL RESULTS

#### 6.3.1 METAWORLD

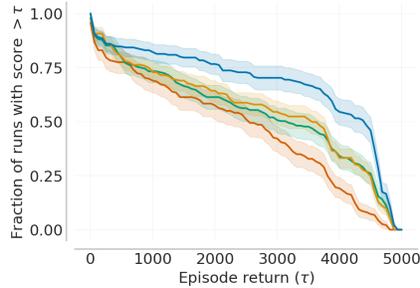


Figure 6: Performance profile of BC (red), CQL (green), CQL+UDS (orange) and CVL (blue) generated by the reliable library ([Agarwal et al., 2021b](#)) for the offline MetaWorld experiments over 5 random seeds.

#### 6.3.2 MOUNTAIN CAR

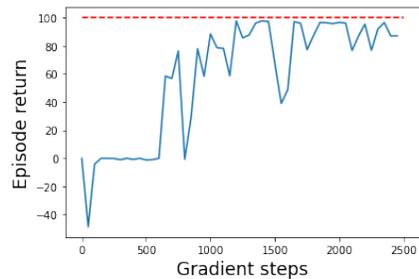


Figure 7: Evaluation returns on Mountain car during training on data from the SAC replay buffer. The red dotted line indicates highest possible return.

Task	BC	CQL	CQL+UDS	CVL (Tabula rasa)	CVL (Pretrained)
basketball	3188 ± 348.9	646 ± 2.2	678 ± 112.6	<b>4503 ± 113.8</b>	4171 ± 285.9
bin-picking	13 ± 1.5	28 ± 0.7	21 ± 1.9	18 ± 1.8	<b>860 ± 69.2</b>
box-close	891 ± 164.4	311 ± 12.4	296 ± 19.8	1496 ± 131.2	<b>4189 ± 352.4</b>
button-press	2667 ± 85.3	3445 ± 60.2	3420 ± 273.4	<b>3659 ± 51.9</b>	1906 ± 360.2
button-press-topdown	3089 ± 256.7	3406 ± 525.7	3505 ± 990.9	<b>3889 ± 36.2</b>	548 ± 49.5
button-press-topdown-wall	1692 ± 47.0	2095 ± 28.4	<b>2135 ± 66.4</b>	2008 ± 21.7	546 ± 90.8
coffee-button	3490 ± 1435.5	3655 ± 740.6	3431 ± 689.8	<b>4259 ± 169.9</b>	149 ± 10.8
coffee-pull	647 ± 11.7	250 ± 10.3	330 ± 3.2	<b>833 ± 27.2</b>	167 ± 0.6
dial-turn	1331 ± 48.7	4257 ± 389.4	4449 ± 276.1	4526 ± 42.8	<b>4611 ± 176.9</b>
disassemble	215 ± 4.3	215 ± 9.6	217 ± 36.0	214 ± 18.6	<b>926 ± 5.6</b>
door-close	3634 ± 141.5	<b>4555 ± 200.9</b>	4547 ± 215.2	4544 ± 7.6	4313 ± 194.0
door-lock	3073 ± 303.7	3775 ± 59.1	<b>3777 ± 144.2</b>	3537 ± 271.1	557 ± 20.6
door-open	828 ± 24.7	4526 ± 71.7	<b>4531 ± 179.0</b>	3985 ± 279.6	613 ± 113.0
door-unlock	1322 ± 181.1	4122 ± 50.2	4002 ± 80.1	3139 ± 413.7	<b>4618 ± 49.7</b>
drawer-close	4619 ± 53.4	4855 ± 0.0	<b>4857 ± 2.0</b>	4853 ± 6.8	2933 ± 671.8
drawer-open	1727 ± 204.0	2768 ± 45.6	2776 ± 25.2	2512 ± 149.3	<b>4664 ± 14.4</b>
faucet-close	4160 ± 49.8	<b>4752 ± 1585.0</b>	4713 ± 1724.2	4683 ± 47.8	4739 ± 57.0
faucet-open	2052 ± 80.9	<b>4731 ± 401.8</b>	4729 ± 561.1	3660 ± 221.9	1637 ± 64.9
hammer	2158 ± 272.0	898 ± 70.3	1030 ± 126.6	<b>4632 ± 73.6</b>	4630 ± 86.5
hand-insert	44 ± 17.8	443 ± 2.0	428 ± 1.5	180 ± 5.3	<b>4612 ± 539.8</b>
handle-press	4734 ± 36.3	2816 ± 4.4	2755 ± 0.8	<b>4861 ± 28.6</b>	2417 ± 169.2
handle-press-side	3820 ± 1556.5	4783 ± 170.5	4786 ± 478.1	<b>4816 ± 352.6</b>	654 ± 27.7
handle-pull	3642 ± 968.8	2422 ± 524.1	2436 ± 1286.8	4594 ± 38.6	<b>4636 ± 35.8</b>
handle-pull-side	3418 ± 1002.2	1898 ± 582.6	1757 ± 343.2	<b>4660 ± 41.0</b>	2904 ± 92.4
lever-pull	3659 ± 180.8	2233 ± 399.5	2157 ± 258.0	<b>4459 ± 107.8</b>	4207 ± 98.9
peg-insert-side	11 ± 1.1	17 ± 4.1	<b>19 ± 1.8</b>	15 ± 0.4	12 ± 0.8
peg-unplug-side	56 ± 1.9	29 ± 2.6	29 ± 2.4	87 ± 1.6	<b>4593 ± 24.6</b>
pick-out-of-hole	10 ± 0.2	207 ± 0.4	191 ± 3.4	<b>1245 ± 186.4</b>	5 ± 0.9
pick-place	1771 ± 416.2	1263 ± 407.6	1306 ± 128.5	2942 ± 454.1	<b>4403 ± 508.3</b>
pick-place-wall	0 ± 0.0	1 ± 0.0	71 ± 0.0	19 ± 0.0	<b>3522 ± 775.3</b>
plate-slide	3979 ± 57.3	2697 ± 475.3	3508 ± 747.0	<b>4649 ± 142.6</b>	802 ± 12.4
plate-slide-back	2402 ± 333.9	3163 ± 1290.3	3014 ± 303.9	<b>4718 ± 306.8</b>	196 ± 5.0
plate-slide-back-side	4017 ± 874.6	4736 ± 1519.0	4732 ± 137.6	<b>4752 ± 196.9</b>	4669 ± 95.1
plate-slide-side	2241 ± 536.9	<b>3104 ± 308.1</b>	3015 ± 329.5	2695 ± 413.8	1939 ± 27.5
push	1834 ± 317.9	494 ± 5.0	463 ± 3.3	1997 ± 196.8	<b>4386 ± 192.7</b>
push-back	9 ± 0.2	71 ± 1.4	135 ± 0.8	109 ± 1.4	<b>204 ± 20.9</b>
push-wall	3327 ± 508.6	689 ± 5.6	628 ± 4.2	4502 ± 176.7	<b>4601 ± 205.6</b>
reach	3069 ± 359.2	3301 ± 920.3	3275 ± 677.8	<b>4819 ± 182.9</b>	4658 ± 204.8
reach-wall	4515 ± 93.9	4828 ± 26.5	<b>4829 ± 49.1</b>	4811 ± 27.0	4825 ± 21.2
stick-pull	595 ± 19.8	297 ± 2.0	441 ± 3.7	<b>4488 ± 52.0</b>	3434 ± 162.9
stick-push	263 ± 6.1	896 ± 3.9	897 ± 3.4	1155 ± 147.5	<b>2804 ± 551.3</b>
sweep	817 ± 124.3	3086 ± 645.7	3162 ± 1507.1	4127 ± 567.2	<b>4461 ± 49.5</b>
sweep-into	532 ± 151.3	1974 ± 34.0	1834 ± 870.1	<b>2657 ± 364.3</b>	506 ± 14.8
window-close	3739 ± 80.4	4478 ± 452.5	4442 ± 13.1	<b>4534 ± 56.2</b>	4519 ± 63.2
window-open	3743 ± 147.3	2773 ± 1433.5	2841 ± 1163.5	<b>4534 ± 109.4</b>	524 ± 320.0

Table 4: Evaluation returns on MetaWorld offline tasks. Average ± standard deviation are shown for 5 random seeds.