

DECOUPLED Q-CHUNKING

Anonymous authors

Paper under double-blind review

ABSTRACT

Bootstrapping bias problem is a long-standing challenge in temporal-difference (TD) methods in off-policy reinforcement learning (RL). Multi-step return backups can alleviate this issue but require delicate importance sampling to correct their off-policy bias. Recent work has proposed to use chunked critics, which estimate the value of short action sequences (“chunks”) rather than individual actions, enabling unbiased multi-step backup. However, extracting policies from chunked critics is challenging: policies must output the entire action chunk open-loop, which can be sub-optimal for environments that require policy reactivity and also challenging to model especially when the chunk length grows. Our key insight is to decouple the chunk length of the critic from that of the policy, allowing the policy to operate over shorter action chunks. We propose a novel algorithm that achieves this by optimizing the policy against a distilled critic for partial action chunks, constructed by optimistically backing up from the original chunked critic to approximate the maximum value achievable when a partial action chunk is extended to a complete one. This design retains the benefits of multi-step value propagation while sidestepping both the open-loop sub-optimality and the difficulty of learning action chunking policies for long action chunks. We evaluate our method on challenging, long-horizon offline goal-conditioned benchmarks and shows that it reliably outperforms prior methods.

1 INTRODUCTION

A reinforcement learning (RL) agent can in principle solve any task with a well-defined reward function, but training an RL agent from scratch can be sample inefficient. In many practical problems, we instead have access to an offline dataset of trajectories that serves as a great prior to accelerate learning. Temporal-difference (TD)-based RL algorithms, which learn a value network to perform approximate dynamic programming via value backups, are particularly suitable in this setting because they are designed to handle off-policy data. A well-known yet long-lasting bottleneck, however, is the bootstrapping bias problem (Jaakkola et al., 1993; Sutton et al., 1998; De Asis et al., 2018; Park et al., 2025)—as the value network regresses towards its own estimates, any error compounds across time steps, making accurate value propagation challenging especially in long-horizon, sparse reward tasks.

Multi-step return backups (such as n -step return (Sutton et al., 1998)) can alleviate bootstrapping bias by effectively reducing the time horizon, but naïvely applying them can result in another form of bias that causes the value estimates to be overly conservative/pessimistic. While it is possible to correct such systematic biases with importance sampling (Munos et al., 2016), they often require additional heuristics and truncations to balance a delicate scale between bias and variance that is often tricky to tune. Recent works (Seo & Abbeel, 2024; Li et al., 2025a; Tian et al., 2025; Li et al., 2025b) leverage chunked value functions, which estimate the value of short action sequences (“chunks”) rather than a single action. This formulation allows n -step return backup without the pessimistic bias (under the open-loop consistency condition, which we will formalize in Section 4). However, directly optimizing a policy over full action chunks is difficult, particularly as the chunk size grows, and it is still unclear how to best extract a policy from a chunked critic.

In this work, we develop a simple, novel technique to address this challenge. We train a policy to predict a shorter, partial action chunk against the chunked critic that takes in longer, complete action chunks. The key design that enables such an optimization is a ‘distilled’ chunked critic with a chunk size that matches the policy: it optimistically regresses to the original chunked critic to approximate the maximum value that the partial action chunk can achieve after being extended into a full action chunk. Conceptually, while the action optimization is still done for the longer, complete action

chunks, the policy network is only trained to output the the partial action chunk of an optimized complete action chunk. This way, the policy only needs to predict a much shorter action chunk (e.g., in the extreme case, only one action), which often admits a much simpler distribution, while enjoying the value learning benefits from the use of chunked critics.

Our main contributions are two-fold. On the theoretical side, we provide a formal analysis of Q-learning with action chunking, identifying the open-loop value learning bias and characterizing the conditions under which action chunking critic backup is preferable over n -step return backup with a single-step critic. On the empirical side, we propose a novel technique, **Decoupled Q-chunking (DQC)**, that addresses the policy learning challenge in action chunking Q-learning by decoupling the policy chunk size from the critic chunk size. DQC trains a policy to only predict a partial action chunk, significantly reducing the policy learning challenge, while retaining the value learning benefits of the chunked critic. We instantiate this technique as a practical offline RL algorithm that outperforms the previous state-of-the-art method on the hardest set of environments in OGBench (Park et al., 2024a), a challenging, long-horizon goal-conditioned RL benchmark.

2 RELATED WORK

Offline and offline-to-online reinforcement learning methods assume access to an offline dataset to learn a policy without interactions with the environment (offline) (Kumar et al., 2020; Kostrikov et al., 2021; Tarasov et al., 2024) or with as little online interaction with the environment as possible (offline-to-online) (Lee et al., 2022; Ball et al., 2023; Nakamoto et al., 2024). Q-learning or TD-based RL algorithms have been a popular choice for these problem settings as they naturally handle off-policy data without the need for on-policy rollouts, and also exhibit great online sample-efficiency (Chen et al., 2021; D’Oro et al., 2022). A large body of literature in these two problem settings has been focusing on tackling the distribution shift challenge by appropriately constraining the policies with respect to the prior offline data, and most of them use the standard 1-step TD backup for Q-learning, which has been known to suffer from the bootstrapping bias problem in the RL literature (Jaakkola et al., 1993; Sutton et al., 1998). To tackle this, recent work (Jeong et al., 2022; Park & Lee, 2024; Park et al., 2025; Li et al., 2025b) has shown that multi-step return backups are effective for improving offline/offline-to-online Q-learning agents. These methods either use a standard single-step critic network (Park et al., 2025) that suffers from the off-policy bias, or use a ‘chunked,’ multi-step critic network (Li et al., 2025b) that does not have such bias but poses a huge policy learning challenge when the chunk size is too large. Our method brings the best of both worlds—it uses action chunking to avoid the off-policy bias while simultaneously avoiding the policy learning challenge by extracting a simpler policy that predicts a shorter action chunk from the full-chunk-sized critic.

Multi-step return backups are computed with multi-step off-policy rewards that can lead to systematic value underestimation (Sutton et al., 1998; Peng & Williams, 1994; Konidaris et al., 2011; Thomas et al., 2015), and there has been a rich literature (Precup et al., 2000; Munos et al., 2016; Rowland et al., 2020) dedicated to fix these biases via importance sampling (Kloek & Van Dijk, 1978) with truncation (Ionides, 2008). These approaches often require a careful balance between bias and variance that can be tricky to tune. More recently, Seo & Abbeel (2024); Li et al. (2025a); Tian et al. (2025); Li et al. (2025b) group temporally extended sequences of actions as chunks and directly estimate the value of an action chunk rather than a single action. Such a formulation allows the value backup to operate directly in the chunk space, which allows multi-step return backup without the systematic biases from the sub-optimal off-policy data. Despite their empirical success, we still lack a good theoretical understanding of the convergence of TD-learning with ‘chunked’ critics, as well as when it should be favored over more traditional multi-step returns. Our work lays out the theoretical foundation for Q-learning with critic chunking, and identifies an important yet subtle, often overlooked bias in the TD-backup. We quantify such bias and provide the condition under which TD backup using critic chunking is guaranteed to perform better than the standard n -step return backup with a single-step critic.

See additional discussions for related work in hierarchical reinforcement learning in Appendix G.

3 PRELIMINARIES

Reinforcement learning can be formalized as a Markov decision process, $\mathcal{M} = (\mathcal{S}, \mathcal{A}, T, r, \rho, \gamma)$, where \mathcal{S} is the state space, \mathcal{A} is the action space, $T : \mathcal{S} \times \mathcal{A} \rightarrow \Delta_{\mathcal{S}}$ is the transition kernel that defines the next state distribution conditioned on the current state and the current action (e.g., $s' \sim T(\cdot | s, a)$), $r : \mathcal{S} \times \mathcal{A} \rightarrow [0, 1]$ is the reward function, $\rho \in \Delta_{\mathcal{S}}$ is the initial state distribution,

and $\gamma \in [0, 1)$ is the discount factor. We also assume we have access to a prior offline dataset $D = \{(s_0^i, a_0^i, r_0^i, s_1^i, a_1^i, r_1^i, \dots, s_H^i)\}_{i=1}^{|D|}$ where the goal is to learn a policy, $\pi : \mathcal{S} \rightarrow \Delta_{\mathcal{A}}$ that maximizes its return, $\eta(\pi) = \mathbb{E}_{s_{t+1} \sim T(\cdot|s_t, a_t), a_t \sim \pi(\cdot|s_t), s_0 \sim \rho} [\sum_{t=0}^{\infty} \gamma^t r(s_t, a_t)]$, the cumulative discounted sum of rewards that the policy receives in expectation.

Temporal difference learning. Modern value-based reinforcement learning methods often learn a critic network, $Q : \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}$ to approximate the maximum discounted cumulative reward starting from state s and action a , and the critic is often trained using the temporal-difference (TD) loss:

$$L(\phi) = \mathbb{E}_{s,a,s' \sim \mathcal{D}} [(Q_\phi(s, a) - r(s, a) - \gamma \bar{Q}(s', a'^*))^2], \quad (1)$$

where \bar{Q} is the target critic that is set to the same critic with its parameters set to an exponential moving average of ϕ , and $a'^* = \arg \max_{a'} Q(s', a')$ (often approximated by a policy π_θ).

Implicit value learning with implicit maximization loss function. Instead of using $Q(s', a'^* \sim \pi_\theta(s'))$ as the TD target, we can use what we refer to as an *implicit maximization* loss function f_{imp} to learn a value function $V_\xi(s)$ that approximates the maximum value $Q(s, a^*)$ (Kostrikov et al., 2021; Hansen-Estruch et al., 2023):

$$L(\xi) = \mathbb{E}_{s,a \sim \mathcal{D}} [f_{\text{imp}}^\kappa(\bar{Q}(s, a) - V_\xi(s))]. \quad (2)$$

Two popular choices of f_{imp}^κ are (1) expectile: $f_{\text{expectile}}^\kappa(c) = |\kappa - \mathbb{I}_{c < 0}|c|^2$, and (2) quantile: $f_{\text{quantile}}^\kappa(c) = |\kappa - \mathbb{I}_{c < 0}||c|$, for any real value $\kappa \in [0.5, 1)$. At the optimum of $L(\xi)$, $V_\xi(s)$ approximates the κ -expectile/quantile of the distribution of the critic values evaluated at $Q(s, a)$, induced by the data distribution \mathcal{D} . With this implicit maximization technique, we no longer need to explicitly find the action a that maximizes $Q(s, a)$ and can use $V_\xi(s)$ as the backup target:

$$L(\phi) = \mathbb{E}_{s,a,s' \sim \mathcal{D}} [(Q_\phi(s, a) - r(s, a) - \gamma V_\xi(s'))^2]. \quad (3)$$

Multi-step return backup. TD learning can sometimes struggle with long-horizon tasks due to the well-known bootstrapping bias problem, where regressing the value network towards its own potentially inaccurate value estimates amplifies the value estimation errors further. To tackle this challenge, we can instead sample a trajectory segment, $(s_t, a_t, s_{t+1}, \dots, a_{t+n-1}, s_{t+n})$, to construct an n -step return backup target from states h steps ahead:

$$L_{\text{ns}}(\phi) = \mathbb{E}_{s_t, a_t, \dots, s_{t+n}} \left[(Q_\phi(s_t, a_t) - R_{t:t+n} - \gamma^n \bar{Q}(s_{t+n}, a_{t+n}^*))^2 \right], \quad (4)$$

where $a_{t+n}^* = \arg \max_{a_{t+n}} Q(s_{t+n}, a_{t+n})$, $R_{t:t+n} := \sum_{t'=t}^{t+n-1} \gamma^{t'-t} r(s_{t'}, a_{t'})$. The n -step return value estimate of reduces the effective horizon by a factor of n , alleviating the bootstrapping bias problem. However, such value estimate is always biased towards the off-policy data distribution, and is also commonly referred to as the *uncorrected n -step return estimator* (Fedus et al., 2020; Kozuno et al., 2021). While there are ways to correct this value estimator via importance sampling (Precup et al., 2000; Munos et al., 2016; Rowland et al., 2020), they require additional tricks (e.g., importance ratio truncation) for numerical stability and re-introduce biases into the estimator, ultimately resulting in a delicate trade-off between variances and biases that must be carefully balanced.

Action chunking critic. Alternatively, one may learn an action chunking critic to estimate the value of a short sequence of actions, $a_{t:t+h} := (a_t, a_{t+1}, \dots, a_{t+h-1})$ (or an *action chunk*) instead: $Q(s_t, a_{t:t+h})$ (Seo & Abbeel, 2024; Li et al., 2025a; Tian et al., 2025; Li et al., 2025b). The TD backup loss for such a critic is naturally multi-step:

$$L_{\text{QC}}(\phi) = \mathbb{E}_{s_t, a_{t:t+h+1}, a_{t:t+h}} \left[(Q_\phi(s_t, a_{t:t+h}) - R_{t:t+h} - \gamma^h \bar{Q}(s_{t+h}, a_{t+h:t+2h}^*))^2 \right], \quad (5)$$

where again $a_{t+h:t+2h}^* = \arg \max_{a_{t+h:t+2h}} Q(s_{t+h}, a_{t+h:t+2h})$. On the one hand, unlike n -step return estimate for single-action critic that is pessimistic, the n -step return estimate (with $n = h$) for the action chunking critic is *unbiased* as long as the action chunk $a_{t:t+h}$ is *independent* of the intermediate states $s_{t+1:t+h+1}$, while enjoying the reduction in effective horizon (Li et al., 2025a;b). On the other hand, action chunking critic implicitly imposes a constraint on the policy that the actions are predicted and executed in chunks. As a result, the policy extracted from the action chunking critic needs to predict the entire action chunk all at once, posing a big learning challenge, especially for environments with complex transition dynamics.

In the following two sections, we offer theoretical insights that characterize the conditions when using action chunking critic is more preferable over n -step return backup with a single critic (Section 4), and develop a practical method that tackles the action chunking policy extraction challenge (Section 5).

4 WHEN SHOULD WE USE ACTION CHUNKING FOR Q-LEARNING?

In this section, we build a theoretical foundation for Q-learning with action chunking critic functions. We start by formalizing the setup of our analysis in Section 4.1, quantifying the value estimation bias incurred from backing up on non-action chunking data (Theorem 4.4) and the optimality of action chunking policy (Theorem 4.6) in Section 4.2. Using these result, we derive the condition when we prefer action chunking Q-learning over the standard n -step return learning in Section 4.3. We also include some examples in which the condition holds in Appendix D in hope to facilitate theoretical analysis of action chunking policy learning in future work.

4.1 ASSUMPTIONS

To build the foundation of our analysis, we start by describing the trajectory data distribution that we use for Q-learning and the trajectory distribution induced by open-loop action chunking policy. In particular, we assume that the trajectory data distribution obeys the transition dynamics T :

Assumption 4.1 (Data Distribution Obeys Dynamics). $\mathcal{D} \in \Delta_{\mathcal{T}}$ is a trajectory distribution generated by rolling out a behavior policy from a distribution of $s_t \sim \mu$. The behavior policy can be non-Markovian (i.e., $\pi_\beta(a_{t+k} \mid s_{t:t+k+1}, a_{t:t+k})$). Each subsequent state is generated obeying the dynamics of the MDP \mathcal{M} : $s_{t+k+1} \sim T(\cdot \mid s_{t+k}, a_{t+k}), \forall k \in \{0, 1, \dots, h-1\}$. The resulting trajectory is $\{s_t, s_{t+1}, \dots, s_{t+h}, a_t, a_{t+1}, \dots, a_{t+h}\} \in \mathcal{T} = \mathcal{S}^h \times \mathcal{A}^h$.

Next, we formally define the open-loop trajectory distribution that we would obtain if we take the same actions in the data and rollout them out open-loop in the environment.

Definition 4.2 (Open-loop Trajectory). From any trajectory distribution \mathcal{D} , we can extract an open-loop policy with a horizon of h by marginalizing out all intermediate states. We use $\pi_{\mathcal{D}}^{\text{open}} : \mathcal{S} \rightarrow \Delta_{\mathcal{A}^h}$ to denote such policy and is formally defined as:

$$\pi_{\mathcal{D}}^{\text{open}}(a_{t:t+h} \mid s_t) := P_{\mathcal{D}}(a_{t:t+h} \mid s_t). \quad (6)$$

By using this open-loop policy to roll-out trajectories in the MDP \mathcal{M} , it induces a trajectory distribution $P_{\mathcal{D}}^{\text{open}} \in \Delta_{\mathcal{S}^{h+1}, \mathcal{A}^h}$ that is generally different from \mathcal{D} . We can decompose this open-loop policy step-by-step with the following factorization $\pi_{\mathcal{D}}^{\text{open}}(a_{t:t+k} \mid s_t) = \prod_{k=0}^{h-1} \pi_{\mathcal{D}}^{\text{open}}(a_{t+k} \mid s_t, a_{t:t+k})$ which allows us to define the induced trajectory distribution $P_{\mathcal{D}}^{\text{open}}$ recursively (for $k \in \{1, 2, \dots, h\}$):

$$P_{\mathcal{D}}^{\text{open}}(s_{t+k}, a_{t:t+k} \mid s_t) := \quad (7)$$

$$P_{\mathcal{D}}^{\text{open}}(s_{t+k-1}, a_{t:t+k-1} \mid s_t) T(s_{t+k} \mid s_{t+k-1}, a_{t+k-1}) \pi_{\mathcal{D}}^{\text{open}}(a_{t+k} \mid s_t, a_{t:t+k}). \quad (8)$$

4.2 OPEN-LOOP VALUE BIAS OF ACTION CHUNKING Q-LEARNING

As what we have elucidated in our definition above, replaying the actions from the trajectory data distribution $P_{\mathcal{D}}$ in an open-loop manner, in general, can result in a different trajectory distribution, $P_{\mathcal{D}}^{\text{open}}$. This discrepancy between $P_{\mathcal{D}}^{\text{open}}$ and $P_{\mathcal{D}}$ has not been carefully analyzed by prior work (e.g., Q-chunking (Li et al., 2025b)) but can play a huge role in the optimal policy that action chunking Q-learning converges to. This is because TD-backup is only unbiased when it is done under the open-loop trajectory distribution $P_{\mathcal{D}}^{\text{open}}$. Naïvely running TD-backup on $P_{\mathcal{D}}$ (as done in Li et al. (2025b)) may lead to a *biased Q-target*. We now formalize the discrepancy and analyze such bias.

Definition 4.3 (Open-Loop Consistency). \mathcal{D} is ε_h -open-loop consistent if for every $s_t \in \mathcal{S}, h' \in \{1, \dots, h\}$, as long as $s_t \in \mathcal{S}$ has non-zero probability in the data (i.e., $P_{\mathcal{D}}(s_t) > 0$),

$$2D_{\text{TV}}(P_{\mathcal{D}}^{\text{open}}(s_{t+h'}, a_{t+h'} \mid s_t) \parallel P_{\mathcal{D}}(s_{t+h'}, a_{t+h'} \mid s_t)) \leq \varepsilon_h, \forall h' \in \{1, 2, \dots, h-1\}, \quad (9)$$

$$2D_{\text{TV}}(P_{\mathcal{D}}^{\text{open}}(s_{t+h} \mid s_t) \parallel P_{\mathcal{D}}(s_{t+h} \mid s_t)) \leq \varepsilon_h. \quad (10)$$

We say \mathcal{D} is strongly ε_h -open-loop consistent if additionally for $h' \in \{1, 2, \dots, h\}$, for every $a_{t:t+h'} \in \mathcal{A}^{h'}$ with non-zero probability in the data (i.e., $P_{\mathcal{D}}(a_{t:t+h'}, s_t) > 0$),

$$2D_{\text{TV}}(T(s_{t+h'} \mid s_t, a_{t:t+h'}) \parallel P_{\mathcal{D}}(s_{t+h'} \mid s_t, a_{t:t+h'})) \leq \varepsilon_h. \quad (11)$$

Intuitively, \mathcal{D} is ε -open-loop consistent if, when executing the same sequence of actions from it open-loop from s_t , the resulting marginal distribution of the state-action h steps into the future (i.e., s_{t+h}) deviates from the corresponding distribution in the dataset by at most ε in total variation distance.

The strong version (Equation (11)) requires the total variation distance bound to hold for every action sequence in the support, whereas the weak version (Equation (9)) only requires the bound to hold in expectation. Having *weak* open-loop consistency of \mathcal{D} is sufficient to show that *behavior* value iteration of an action chunking critic results in a *nominal* value function with a bounded bias from the true value of the open-loop policy $\pi_{\mathcal{D}}^{\text{open}}$:

Theorem 4.4 (Bias of Action Chunking Critic). *Let $\hat{V}_{\text{ac}} : \mathcal{S} \rightarrow [0, 1/(1 - \gamma)]$ be a solution of*

$$\hat{V}_{\text{ac}}(s_t) = \mathbb{E}_{s_{t+1:t+h+1}, a_{t:t+h} \sim P_{\mathcal{D}}(\cdot | s_t)} \left[R_{t:t+h} + \gamma^h \hat{V}_{\text{ac}}(s_{t+h}) \right], \quad (12)$$

with $R_{t:t+h} = \sum_{t'=t}^{t+h} \gamma^{t'-t} r(s_{t'}, a_{t'})$ and V_{ac} is the true value of $\pi_{\mathcal{D}}^{\text{open}} : s_t \mapsto P_{\mathcal{D}}(a_{t:t+h} | s_t)$. If \mathcal{D} is ε_h -open-loop consistent, then under $\text{supp}(\mathcal{D})$,

$$\|V_{\text{ac}} - \hat{V}_{\text{ac}}\|_{\infty} \leq \frac{\varepsilon_h}{(1 - \gamma^h)(1 - \gamma)}. \quad (13)$$

The proof of Theorem 4.4 is available in Appendix E. A direct consequence of this result is that the true value of the optimal action chunking policy is close to that of the optimal closed-loop policy:

Corollary 4.5 (Optimal Action Chunking Policy). *Let $\pi^* : \mathcal{S} \rightarrow \Delta_{\mathcal{A}}$ be an optimal policy in \mathcal{M} and \mathcal{D}^* be the data collected by π^* . If \mathcal{D}^* is ε_h -open-loop consistent, then under $\text{supp}(\mathcal{D}^*)$,*

$$\|V_{\text{ac}}^* - V^*\|_{\infty} \leq \|\tilde{V}_{\text{ac}} - V^*\|_{\infty} \leq \frac{\varepsilon_h}{(1 - \gamma^h)(1 - \gamma)}, \quad (14)$$

where V^ is the value of the optimal policy π^* , V_{ac}^* is the true value of the optimal action chunking policy, and \tilde{V}_{ac} is the true value of the action chunking policy from cloning the data \mathcal{D}^* :*

$$\tilde{\pi}_{\text{ac}}(a_{t:t+h} | s_t) : s_t \mapsto P_{\mathcal{D}^*}(\cdot | s_t). \quad (15)$$

The proof of Corollary 4.5 (available in Appendix E) builds on the observation that the nominal (biased) value of the action chunking critic obtained from behavior value iteration on an optimal data \mathcal{D}^* (i.e., the data collected from an optimal policy π^*) recovers the value of the optimal policy. This allows us to use Theorem 4.4 to show that the value of the action chunking policy obtained by behavior cloning on such optimal data is close to the nominal (biased) value of its critic, and thus close to the optimal value of the closed-loop policy.

Next, we analyze the performance of the action chunking policy obtained by Q-learning. In particular, we analyze the Q-function obtained as a solution of the following equation under $\text{supp}(\mathcal{D})$:

$$\hat{Q}_{\text{ac}}^+(s_t, a_{t:t+h}) = \mathbb{E}_{s_{t+1:t+h+1} \sim P_{\mathcal{D}}(\cdot | s_t, a_{t:t+h})} \left[R_{t:t+h} + \gamma^h \max_{a_{t+h:t+2h}} \hat{Q}_{\text{ac}}^+(s_{t+h}, a_{t+h:t+2h}) \right]. \quad (16)$$

The corresponding action chunking policy is

$$\pi_{\text{ac}}^+ : s_t \mapsto \arg \max_{a_{t:t+h}} \hat{Q}_{\text{ac}}^+(s_t, a_{t:t+h}). \quad (17)$$

It turns out that with the weak version of the open-loop consistent condition, the worst case performance of the action chunking policy may be arbitrarily low (see an example in Appendix F). Fortunately, as long as the data \mathcal{D} satisfies the strongly open-loop consistency (Equation (11)), we can show that the learned policy π_{ac}^+ is provably near-optimal by combining all the results above together:

Theorem 4.6 (Q-Learning with Action Chunking Policy on Off-policy Data). *If \mathcal{D} is strongly ε_h -open-loop consistent and $\text{supp}(\mathcal{D}) \supseteq \text{supp}(\mathcal{D}^*)$, with \mathcal{D}^* being the data distribution of an arbitrary optimal policy π^* under \mathcal{M} , then the following bound holds under $\text{supp}(\mathcal{D}^*)$:*

$$\|V_{\text{ac}}^+ - V^*\|_{\infty} \leq \frac{3\varepsilon_h}{(1 - \gamma^h)(1 - \gamma)}, \quad (18)$$

where V^ is the value of an optimal policy under \mathcal{M} .*

The implication of Theorem 4.6 (proof available in Appendix E) is that as long as \mathcal{D} satisfies the strongly open-loop consistency condition and contains the behavior in \mathcal{D}^* , Q-learning with action chunking is guaranteed to converge to a near-optimal action chunking policy regardless of how sub-optimal the data \mathcal{D} might be. As we will show in the following section, this is in contrast to n -step return policy where its performance depends on the sub-optimality of the data.

4.3 COMPARING TO n -STEP RETURN Q-LEARNING

We now characterize the condition when action chunking Q-learning should be preferred over the standard n -step return backup. We start by introducing a notion of sub-optimality of the data \mathcal{D} :

Definition 4.7 (Sub-optimal data). \mathcal{D} is δ_n -suboptimal for backup horizon length $n \in \mathbb{N}^+$ if

$$Q^*(s_t, a_t) - \mathbb{E}_{P_{\mathcal{D}}(\cdot|s_t, a_t)} [R_{t:t+n} + \gamma^n V^*(s_{t+n})] \geq \delta_n, \forall s_t \in \mathcal{S}, a_t \in \mathcal{A}. \quad (19)$$

Intuitively, δ_n captures how much worse the n -step return policy can get compared to the optimal policy incurred by the backup bias. Under such condition, we can show that the action chunking policy is provably better than the n -step return policy as long as δ_n is large.

Theorem 4.8. Let \mathcal{D} be strongly ε_h -open-consistent, δ_n -suboptimal, and $\text{supp}(\mathcal{D}) \supseteq \text{supp}(\mathcal{D}^*)$. Let π_n^* be the optimal n -step return policy learned from \mathcal{D} , as the solution of

$$Q_n^*(s_t, a_t) = \mathbb{E}_{P_{\mathcal{D}}} [R_{t:t+n} + \gamma^n Q_n^*(s_{t+n}, \pi_n^*(s_{t+n}))], \quad \pi_n^*: s_t \mapsto \arg \max_{a_t} Q_n^*(s_t, a_t). \quad (20)$$

As long as $\delta_n > \frac{3\varepsilon_h(1-\gamma^n)}{(1-\gamma)(1-\gamma^h)}$, then from all $s \in \text{supp}(\mathcal{D}^*)$, the action chunking policy, π_{ac}^+ (Equation (17)), is better than the n -step return policy, π_n (Equation (20)) (i.e., $V_{ac}^+(s) > V_n^*(s)$).

The proof of Theorem 4.8 is available in Appendix E. Notably, for $n = h$, the condition on δ_n and ε_h reduces to $\delta_n > 3\varepsilon_h H$ with effective horizon H (i.e., $H = 1/(1-\gamma)$). As long as \mathcal{D} is more than $O(\varepsilon_h H)$ sub-optimal, the action chunking policy performs provably better than n -step return policy.

5 DECOUPLED Q-CHUNKING

We propose a new algorithm that enjoys the benefits of value backup speedup of Q-chunking while avoiding the difficulty of learning an open-loop action chunking policy with a large chunk size.

Our core idea is to decouple the chunk size of the critic from that of the policy. In particular, we train a policy $\pi(a_{t:t+h_a} | s_t)$ to output an action chunk (with a size of $h_a \ll h$) with the following objective:

$$L(\pi) := -\mathbb{E}_{a_{t:t+h_a} \sim \pi(\cdot|s_t)} [Q_\phi(s, [a_{t:t+h_a}, a_{t+h_a:t+h}^*])], \quad (21)$$

where $[a_{t:t+h_a}, a_{t+h_a:t+h}^*]$ represents the concatenation of two partial action chunks (size h_a and size $h - h_a$) into a full action chunk $a_{t:t+h}$ of size h , and $a_{t+h_a:t+h}^*$ is the best ‘second-half’ of the action chunk that maximizes the critic value under Q_ϕ :

$$a_{t+h_a:t+h}^* := \arg \max_{a_{t+h_a:t+h}} Q_\phi(s, [a_{t:t+h_a}, a_{t+h_a:t+h}]). \quad (22)$$

Essentially, we want our policy to predict the partial action chunk (of size h_a) within an optimal action chunk of size h , rather than the entire optimal action chunk. This lowers the policy expressivity requirement and hence the learning challenges associated with it with $h_a < h$.

However, directly optimizing this objective (Equation (21)) does not lead to a novel algorithm because taking the maximization over $a_{t+h_a:t+h}$ seemingly requires us to learn a policy of the original chunk size anyways. To address this issue, we learn a separate partial critic Q_ψ^P , which only takes in the partial action chunk (of size h_a) as input, to approximate the maximum value this partial action chunk can achieve when it is extended to the full action chunk (of size h):

$$Q_\psi^P(s, a_{t:t+h_a}) \approx Q_\phi(s, [a_{t:t+h_a}, a_{t+h_a:t+h}^*]) \quad (23)$$

To train Q_ψ^P , we can use an *implicit maximization* loss function (as described in Equation (2)):

$$L(\psi) := f_{\text{imp}}^{\kappa_d}(\bar{Q}_\phi(s_t, a_{t:t+h}) - Q_\psi^P(s_t, a_{t:t+h_a})), \quad (24)$$

where $s_t, a_{t:t+h}$ are sampled from \mathcal{D} . As a result, the partial critic, Q_ψ^P , is distilled from the original critic via an optimistic regression, where its optimum $Q_\psi^*(s, a_{t:t+h_a})$ approximates $Q_\phi(s, [a_{t:t+h_a}, a_{t+h_a:t+h}^*])$ in Equation (21), conveniently removing the need for training a policy to predict the whole optimal action chunk entirely. This allows us to simplify the policy objective as

$$L(\pi) := -\mathbb{E}_{a_{t:t+h_a} \sim \pi(\cdot|s_t)} [Q_\psi^P(s, a_{t:t+h_a})]. \quad (25)$$

In summary, DQC trains a policy to predict a partial chunk, $a_{t:t+h_a}$ (of size h_a), by hill climbing the value of a partial critic $Q_\psi^P(s, a_{t:t+h_a})$ that is distilled from the original chunked critic $Q_\phi(s, a_{t:t+h})$

Algorithm 1 Decoupled Q-chunking (DQC).

Given: $D, Q_\phi(s_t, a_{t:t+h}), Q_\psi^P(s_t, a_{t:t+h_a}), V_\xi(s_t), \pi_\beta(a_{t:t+h_a} | s_t)$

1. Agent Update:
 $(s_{t:t+h+1}, a_{t:t+h}, r_{t:t+h}) \sim D.$ \triangleright sample trajectory chunk from the offline dataset
 Optimize Q_ϕ with $L(\phi) = \left(Q_\phi(s_t, a_{t:t+h}) - \sum_{k=0}^{h-1} \gamma^k r_{t+k} - \gamma^h \bar{V}_\xi(s_{t+h}) \right)^2.$
 Optimize Q_ψ^P with $L(\psi) = f_{\text{expectile}}^{\tau_d}(\bar{Q}_\phi(s_t, a_{t:t+h}) - Q_\psi^P(s_t, a_{t:t+h_a})).$
 Optimize V_ξ with $L(\xi) = f_{\text{quantile}}^{\tau_b}(\bar{Q}_\psi^P(s_t, a_{t:t+h_a}^\beta) - V_\xi(s_t)), a_{t:t+h_a}^\beta \sim \pi_\beta(\cdot | s_t)$

2. Policy Extration:
 $a_{t:t+h_a}^1, a_{t:t+h_a}^2, \dots, a_{t:t+h_a}^N \sim \pi_\beta(\cdot | s_t)$ \triangleright sample N actions from behavior policy
 $a_{t:t+h_a}^* \leftarrow \arg \max_{\{a_{t:t+h_a}^i\}_{i=1}^N} Q_\psi^P(s_t, a_{t:t+h_a}^i)$ \triangleright take the action with the highest Q-value

via an implicit maximization loss. This allows our policy to fully leverage the chunked critic Q_ϕ (and thus the value speedup benefits associated with Q-chunking) without the need to predict the full action chunk (of size h), mitigating the learning challenge of an action chunking policy.

Practical considerations for offline RL. Finally, we describe several implementation details that we find to work well in the offline RL setting, which our experiments primarily focus on. Our implementation draws inspirations from a prior method, IDQL (Hansen-Estruch et al., 2023).

We first train a behavior cloning flow policy π_β using a standard flow-matching objective (Liu et al., 2022) on the offline dataset D . Then, we approximate the policy optimization objective in DQC (Equation (25)) using best-of- N sampling without explicitly modeling π :

$$a_{t:t+h_a}^* \leftarrow \arg \max_{\{a_{t:t+h_a}^i\}_{i=1}^N} Q_\psi^P(s_t, a_{t:t+h_a}^i), \quad \text{where } a_{t:t+h_a}^1, \dots, a_{t:t+h_a}^N \sim \pi_\beta(\cdot | s_t). \quad (26)$$

where $a_{t:t+h_a}^*$ is output of the policy that we extract from Q_ψ^P for state s_t . Essentially, this sampling procedure is a test-time approximation of the objective in Equation (25), where it outputs action (chunk) that maximizes Q_ψ^P , subject to the behavior prior, as modeled by π_β .

For TD learning of Q_ϕ , directly computing the TD backup target from either Q_ϕ or Q_ψ^P is computationally expensive, as either requires samples from the current policy, which is approximated via the best-of- N sampling procedure as described above. Instead, we use the implicit value backup (Kostrikov et al., 2021) (i.e., as described in Equation (2)) to approximate the target:

$$L(\xi) = f_{\text{quantile}}^{\kappa_b}(\bar{Q}_\psi^P(s_t, a_{t:t+h_a}^\beta) - V_\xi(s_t)), \quad a_{t:t+h_a}^\beta \sim \pi_\beta(\cdot | s_t) \quad (27)$$

where we pick the quantile regression loss as the implicit maximization loss function. This is because the Q-value obtained from best-of- N sampling can be seen as the largest order statistic of a random batch (of size N) of the behavior Q-values (i.e., $\{Q(s, a^i)\}_{i=1}^N, a^i \sim \pi_\beta(\cdot | s)$). Such statistic estimates the behavior Q-value distribution’s $\frac{N}{1-N}$ -quantile, which is the same as $V_\xi(s)$ at the optimum of $L(\xi)$ if we set $\kappa_b = \frac{N}{1-N}$. In practice, we use a larger κ_b for numerical stability.

Finally, we pick the expectile regression loss for training the distilled partial critic Q_ψ^P because prior work has found it to work the best among all implicit maximization loss functions (Hansen-Estruch et al., 2023). A summary of the algorithm is available in Algorithm 1.

6 EXPERIMENTAL SETUP

We conduct experiments to evaluate the benefits of decoupling the policy chunk size and the critic chunk size on OGBench (Park et al., 2024a)—a challenging long-horizon, goal-conditioned offline RL benchmark consisting of a diverse set of environments (from manipulation to locomotion). In particular, we use the more difficult environments introduced by Park et al. (2025) (Figure 4), where multi-step return backups are crucial. These environments require highly complex, long-horizon reasoning. For example, the puzzle tasks require stitching up to 24 atomic motions to solve a combinatorial puzzle with a robot arm, and the humanoidmaze task requires controlling a high-dimensional humanoid robot over 3000 environment steps to navigate a maze. These environments serve as an ideal testbed for our algorithm, which improves upon n -step returns and Q-chunking. We now describe our main comparisons. To start with, we consider several direct ablation baselines:

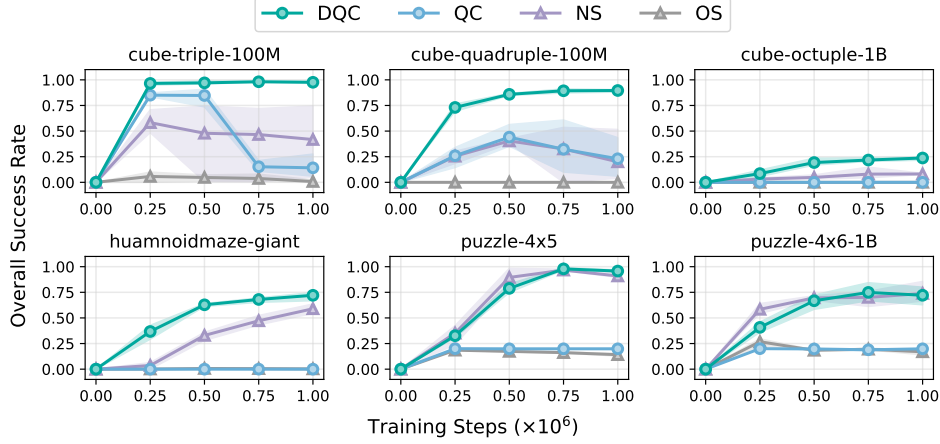


Figure 1: **Offline goal-conditioned RL results.** Our method (*DQC*) uses a *decoupled* critic and policy chunk sizes, which allows it to consistently outperform our baselines: *QC*: Q-chunking (Li et al., 2025b); *NS*: *n*-step return backup; *OS*: 1-step TD-backup.

QC (Li et al., 2025b) uses a single critic that has the same chunk length as that of the policy (*i.e.*, $h = h_a$). This baseline tests whether having *decoupled* chunk sizes is important.

NS: *n*-step return TD backup. This baseline uses a single one-step critic (*i.e.*, $Q(s_t, a_t)$). Compared to *DQC* with $h = n$ and $h_a = 1$, this baseline tests whether using a chunked critic is important.

OS: Standard 1-step TD backup. This is the same as *NS* but with $n = 1$.

Beyond the ablation baselines, we also consider the following strong goal-conditioned baselines:

FBC/HFBC: Goal-conditioned and hierarchical goal-conditioned flow behavior cloning baselines considered in Park et al. (2025).

IQL/HIQL (Kostrikov et al., 2021; Park et al., 2023): These are strong goal-conditioned RL methods that train a goal-conditioned value function with implicit value backups and extract a flat (IQL) or hierarchical (HIQL) policy from the value function.

SHARSA (Park et al., 2025): The previous state-of-the-art method on the long-horizon environments that we evaluate on. The method uses a combination of *n*-step return and bi-level hierarchical policies.

In our ablation study, we also consider an additional baseline, **QC-NS**, that uses the idea of decoupled policy chunking and critic chunking ($h_a < h$), but without using a distilled critic. This baseline simply uses *n*-step return targets to directly train a critic with a chunk size of h_a without implicit maximization (Equation (24)). The performance of this baseline helps determine how important it is to learn a separate distilled critic for partial action chunks with implicit maximization. For all our main results, we run 3 seeds and report the means and the 95% confidence intervals.

Task	FBC	HFBC	IQL	HIQL	SHARSA	OS	NS	QC	DQC
cube-triple-100M	53 _[48,57]	57 _[54,61]	64 _[59,68]	-	82 _[78,88]	1 _[0,2]	42 _[5,74]	14 _[6,28]	98 _[96,98]
cube-quadruple-100M	32 _[30,33]	38 _[34,41]	53 _[53,53]	-	67 _[62,74]	0 _[0,0]	20 _[4,52]	23 _[6,44]	90 _[88,90]
cube-octuple-1B	0 _[0,0]	20 _[17,23]	0 _[0,0]	1 _[0,2]	20 _[19,20]	0 _[0,0]	8 _[8,9]	0 _[0,0]	24 _[22,25]
humanoidmaze-giant	1 _[0,2]	19 _[16,22]	3 _[3,4]	22 _[18,29]	18 _[13,25]	0 _[0,1]	59 _[54,64]	0 _[0,0]	72 _[67,75]
puzzle-4x5	0 _[0,0]	4 _[2,6]	19 _[18,20]	5 _[3,7]	1 _[0,2]	14 _[13,15]	91 _[90,94]	20 _[19,20]	96 _[95,97]
puzzle-4x6-1B	0 _[0,0]	2 _[1,2]	17 _[16,19]	8 _[4,12]	56 _[49,63]	17 _[14,19]	74 _[66,86]	20 _[20,20]	72 _[63,80]

Table 1: **Comparisons with prior methods.** Our method outperforms SHARSA (the previous state-of-the-art method on this benchmark) on all environments. (*) indicates that we take the results from the original paper (Park et al., 2025), where we take the results with larger 10M-sized datasets for humanoidmaze-giant (originally 4M) and puzzle-4x5 (originally 3M). We omit HIQL results on cube-{triple, quadruple} given its high computational cost and poor performance on the other tasks.

7 RESULTS

In this section, we present our experimental results to answer the following three questions:

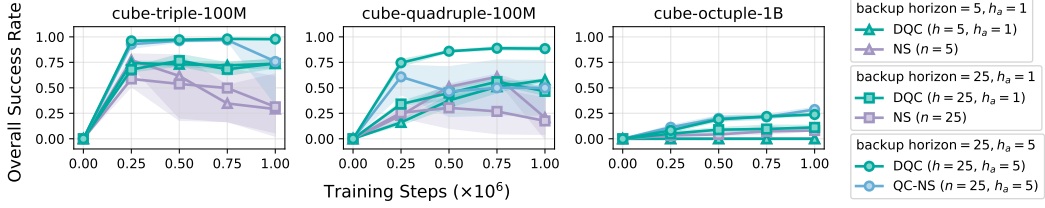


Figure 2: **Distilled critic ablations.** Each group in the legend contains DQC and its non-distilled counterpart with the same configuration (*i.e.*, same backup horizon and same policy chunk size). Our method (DQC) performs on par or better than the non-distilled counterpart across all configurations.

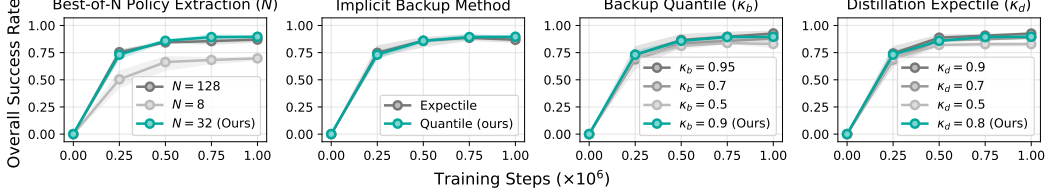


Figure 3: **Hyperparameter sensitivity analysis on cube-quadruple-100M.** *Best-of-N* (N): the number of action samples drawn from $\pi_\beta(\cdot | s)$ during policy evaluation; *Implicit Backup Method*: the implicit maximization loss function for the implicit value backup; *Backup Quantile* (κ_b): the coefficient for the implicit value backup; *Distillation Expectile* (κ_d): the coefficient for training the distilled critic.

(Q1) Does DQC improve upon n -step return, Q-chunking? Figure 1 compares DQC (ours) to both n -step and QC across six challenging long-horizon GCRL tasks, with our method performing on par or better across the board. Table 1 shows DQC also consistently outperforms the previous state-of-the-art method on this benchmark, SHARSA (Park et al., 2025), on all environments. For each environment, we pick the best configuration (in terms of h , h_a , and n). See Appendix C (DQC: Table 6, SHARSA/QC/NS: Table 7) for the environment-specific hyperparameters used in Figure 1 and Table 1. For the results on all configurations, see the complete table in Appendix A.

(Q2) Is training a separate distilled critic Q_ψ^P necessary? In Figure 2, we compare DQC without using the distilled critic across three different (h, h_a) configurations: $(h = 25, h_a = 5)$, $(h = 25, h_a = 1)$, and $(h = 5, h_a = 1)$. For configurations with $h_a = 1$, the baseline without using the distilled critic is the same as the n -step return baseline (with $n = h$) and for the configuration with $h_a = 5$, it is the same as combining Q-chunking and n -step return. Across three configurations, DQC performs on par or better than its non-distilled counterpart. This highlights that the use of a separate distilled critic for the partial action chunk is necessary for the effectiveness of DQC.

(Q3) How sensitive is DQC to its hyperparameters? Figure 3 shows that our method is neither sensitive to the implicit backup method (quantile or expectile), nor sensitive to the backup coefficient κ_b . The important hyperparameters are the N in best-of-N policy extraction and the distillation expectile coefficient, κ_d . Making sure the number of action samples N is large enough (*e.g.*, $N = 32$) is crucial for good performance, though a larger N ($N = 128$) does not lead to better performance.

8 DISCUSSION

We provide a theoretical foundation for action chunking Q-learning and demonstrate how to effectively extract policies from chunked critics. Theoretically, we provide a formal analysis of action chunking Q-learning, identifying the TD backup bias that arises from *open-loop inconsistency* and characterizing the conditions under which action chunking Q-learning is preferred over n -step return learning. Empirically, we develop a novel technique that enables effective policy extraction from chunked critics with long action chunks, *scaling up action chunking Q-learning* to much harder environments. Together, these contributions advance the goal of tackling bootstrapping bias in TD-learning. Several challenges remain, indicating promising avenues for future research. Our method still inherits the open-loop value bias identified in Theorem 4.4, and developing techniques to actively correct for this bias could further improve performance. Moreover, our method relies on a fixed policy action chunk size h_a and critic action chunk size h across all states, even though the optimal action chunk size may vary by state. Developing practical methods that can support flexible, state-dependent chunk sizes would be a natural next step.

REPRODUCIBILITY STATEMENT

To facilitate future research, we include our source code as part of the supplementary materials, along with example scripts for both our method and our baselines. We describe our environments in Appendix B and hyperparameters in Appendix C. For our theoretical results, we fully state our assumption in Assumption 4.1 and provide complete proofs in Appendix E.

REFERENCES

- Anurag Ajay, Aviral Kumar, Pulkit Agrawal, Sergey Levine, and Ofir Nachum. OPAL: Offline primitive discovery for accelerating offline reinforcement learning. In *International Conference on Learning Representations*, 2021. URL <https://openreview.net/forum?id=V69LGwJ01IN>.
- Pierre-Luc Bacon, Jean Harb, and Doina Precup. The option-critic architecture. In *Proceedings of the AAAI conference on artificial intelligence*, volume 31, 2017.
- Akhil Bagaria and George Konidaris. Option discovery using deep skill chaining. In *International Conference on Learning Representations*, 2019.
- Akhil Bagaria, Ben Abbatematto, Omer Gottesman, Matt Corsaro, Sreehari Rammohan, and George Konidaris. Effectively learning initiation sets in hierarchical reinforcement learning. *Advances in Neural Information Processing Systems*, 36, 2024.
- 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.
- Boyuan Chen, Chuning Zhu, Pulkit Agrawal, Kaiqing Zhang, and Abhishek Gupta. Self-supervised reinforcement learning that transfers using random features. *Advances in Neural Information Processing Systems*, 36, 2024.
- Xinyue Chen, Che Wang, Zijian Zhou, and Keith Ross. Randomized ensembled double q-learning: Learning fast without a model. *arXiv preprint arXiv:2101.05982*, 2021.
- Nuttapong Chentanez, Andrew Barto, and Satinder Singh. Intrinsically motivated reinforcement learning. *Advances in neural information processing systems*, 17, 2004.
- Christian Daniel, Gerhard Neumann, Oliver Kroemer, and Jan Peters. Hierarchical relative entropy policy search. *Journal of Machine Learning Research*, 17(93):1–50, 2016.
- Peter Dayan and Geoffrey E Hinton. Feudal reinforcement learning. *Advances in neural information processing systems*, 5, 1992.
- Kristopher De Asis, J Hernandez-Garcia, G Holland, and Richard Sutton. Multi-step reinforcement learning: A unifying algorithm. In *Proceedings of the AAAI conference on artificial intelligence*, volume 32, 2018.
- Anita de Mello Koch, Akhil Bagaria, Bingnan Huo, Zhiyuan Zhou, Cameron Allen, and George Konidaris. Learning transferable sub-goals by hypothesizing generalizing features. 2025.
- Thomas G Dietterich. Hierarchical reinforcement learning with the maxq value function decomposition. *Journal of artificial intelligence research*, 13:227–303, 2000.
- 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 *Deep Reinforcement Learning Workshop NeurIPS 2022*, 2022.
- Ishan P Durugkar, Clemens Rosenbaum, Stefan Dernbach, and Sridhar Mahadevan. Deep reinforcement learning with macro-actions. *arXiv preprint arXiv:1606.04615*, 2016.
- William Fedus, Prajit Ramachandran, Rishabh Agarwal, Yoshua Bengio, Hugo Larochelle, Mark Rowland, and Will Dabney. Revisiting fundamentals of experience replay. In *International conference on machine learning*, pp. 3061–3071. PMLR, 2020.

- Roy Fox, Sanjay Krishnan, Ion Stoica, and Ken Goldberg. Multi-level discovery of deep options. *arXiv preprint arXiv:1703.08294*, 2017.
- Kevin Frans, Seohong Park, Pieter Abbeel, and Sergey Levine. Unsupervised zero-shot reinforcement learning via functional reward encodings. In Ruslan Salakhutdinov, Zico Kolter, Katherine Heller, Adrian Weller, Nuria Oliver, Jonathan Scarlett, and Felix Berkenkamp (eds.), *Proceedings of the 41st International Conference on Machine Learning*, volume 235 of *Proceedings of Machine Learning Research*, pp. 13927–13942. PMLR, 21–27 Jul 2024. URL <https://proceedings.mlr.press/v235/frans24a.html>.
- Jonas Gehring, Gabriel Synnaeve, Andreas Krause, and Nicolas Usunier. Hierarchical skills for efficient exploration. *Advances in Neural Information Processing Systems*, 34:11553–11564, 2021.
- Philippe Hansen-Estruch, Ilya Kostrikov, Michael Janner, Jakub Grudzien Kuba, and Sergey Levine. Idql: Implicit q-learning as an actor-critic method with diffusion policies. *arXiv preprint arXiv:2304.10573*, 2023.
- Hao Hu, Yiqin Yang, Jianing Ye, Ziqing Mai, and Chongjie Zhang. Unsupervised behavior extraction via random intent priors. In *Thirty-seventh Conference on Neural Information Processing Systems*, 2023. URL <https://openreview.net/forum?id=4vGVQVz5KG>.
- Edward L Ionides. Truncated importance sampling. *Journal of Computational and Graphical Statistics*, 17(2):295–311, 2008.
- Tommi Jaakkola, Michael Jordan, and Satinder Singh. Convergence of stochastic iterative dynamic programming algorithms. *Advances in neural information processing systems*, 6, 1993.
- Jihwan Jeong, Xiaoyu Wang, Michael Gimelfarb, Hyunwoo Kim, Baher Abdulhai, and Scott Sanner. Conservative bayesian model-based value expansion for offline policy optimization. *arXiv preprint arXiv:2210.03802*, 2022.
- Teun Kloek and Herman K Van Dijk. Bayesian estimates of equation system parameters: an application of integration by monte carlo. *Econometrica: Journal of the Econometric Society*, pp. 1–19, 1978.
- George Konidaris, Scott Niekum, and Philip S Thomas. TD_{γ} : Re-evaluating complex backups in temporal difference learning. *Advances in Neural Information Processing Systems*, 24, 2011.
- George Dimitri Konidaris. *Autonomous robot skill acquisition*. University of Massachusetts Amherst, 2011.
- Ilya Kostrikov, Ashvin Nair, and Sergey Levine. Offline reinforcement learning with implicit Q-learning. *arXiv preprint arXiv:2110.06169*, 2021.
- Tadashi Kozuno, Yunhao Tang, Mark Rowland, Rémi Munos, Steven Kapturowski, Will Dabney, Michal Valko, and David Abel. Revisiting Peng’s Q (λ) for modern reinforcement learning. In *International Conference on Machine Learning*, pp. 5794–5804. PMLR, 2021.
- Tejas D Kulkarni, Karthik Narasimhan, Ardavan Saeedi, and Josh Tenenbaum. Hierarchical deep reinforcement learning: Integrating temporal abstraction and intrinsic motivation. *Advances in neural information processing systems*, 29, 2016.
- Aviral Kumar, Aurick Zhou, George Tucker, and Sergey Levine. Conservative Q-learning for offline reinforcement learning. *Advances in Neural Information Processing Systems*, 33:1179–1191, 2020.
- Seunghyun Lee, Younggyo Seo, Kimin Lee, Pieter Abbeel, and Jinwoo Shin. Offline-to-online reinforcement learning via balanced replay and pessimistic Q-ensemble. In *Conference on Robot Learning*, pp. 1702–1712. PMLR, 2022.
- Ge Li, Dong Tian, Hongyi Zhou, Xinkai Jiang, Rudolf Lioutikov, and Gerhard Neumann. TOP-ERL: Transformer-based off-policy episodic reinforcement learning. In *The Thirteenth International Conference on Learning Representations*, 2025a. URL <https://openreview.net/forum?id=N4NhVN30ph>.

- Qiyang Li, Zhiyuan Zhou, and Sergey Levine. Reinforcement learning with action chunking. *arXiv preprint arXiv:2507.07969*, 2025b.
- Xingchao Liu, Chengyue Gong, and Qiang Liu. Flow straight and fast: Learning to generate and transfer data with rectified flow. *arXiv preprint arXiv:2209.03003*, 2022.
- Amy McGovern and Richard S Sutton. Macro-actions in reinforcement learning: An empirical analysis. 1998.
- Ishai Menache, Shie Mannor, and Nahum Shimkin. Q-cut—dynamic discovery of sub-goals in reinforcement learning. In *Machine Learning: ECML 2002: 13th European Conference on Machine Learning Helsinki, Finland, August 19–23, 2002 Proceedings 13*, pp. 295–306. Springer, 2002.
- Josh Merel, Leonard Hasenclever, Alexandre Galashov, Arun Ahuja, Vu Pham, Greg Wayne, Yee Whye Teh, and Nicolas Heess. Neural probabilistic motor primitives for humanoid control. *arXiv preprint arXiv:1811.11711*, 2018.
- Rémi Munos, Tom Stepleton, Anna Harutyunyan, and Marc Bellemare. Safe and efficient off-policy reinforcement learning. *Advances in neural information processing systems*, 29, 2016.
- Ofir Nachum, Shixiang Shane Gu, Honglak Lee, and Sergey Levine. Data-efficient hierarchical reinforcement learning. *Advances in neural information processing systems*, 31, 2018.
- 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. *Advances in Neural Information Processing Systems*, 36, 2024.
- Soroush Nasiriany, Tian Gao, Ajay Mandlekar, and Yuke Zhu. Learning and retrieval from prior data for skill-based imitation learning. In *Conference on Robot Learning*, 2022.
- Alexandros Paraschos, Christian Daniel, Jan R Peters, and Gerhard Neumann. Probabilistic movement primitives. *Advances in neural information processing systems*, 26, 2013.
- Kwanyoung Park and Youngwoon Lee. Model-based offline reinforcement learning with lower expected q-learning. *arXiv preprint arXiv:2407.00699*, 2024.
- Seohong Park, Dibya Ghosh, Benjamin Eysenbach, and Sergey Levine. HIQL: Offline goal-conditioned RL with latent states as actions. In *Thirty-seventh Conference on Neural Information Processing Systems*, 2023. URL <https://openreview.net/forum?id=cLQCCTVDuW>.
- Seohong Park, Kevin Frans, Benjamin Eysenbach, and Sergey Levine. Ogbench: Benchmarking offline goal-conditioned rl. *ArXiv*, 2024a.
- Seohong Park, Tobias Kreiman, and Sergey Levine. Foundation policies with hilbert representations. In *Forty-first International Conference on Machine Learning*, 2024b. URL <https://openreview.net/forum?id=LhNsSaAKub>.
- Seohong Park, Kevin Frans, Deepinder Mann, Benjamin Eysenbach, Aviral Kumar, and Sergey Levine. Horizon reduction makes RL scalable. *arXiv preprint arXiv:2506.04168*, 2025.
- Jing Peng and Ronald J Williams. Incremental multi-step Q-learning. In *Machine Learning Proceedings 1994*, pp. 226–232. Elsevier, 1994.
- Xue Bin Peng, Glen Berseth, KangKang Yin, and Michiel Van De Panne. Deeploco: Dynamic locomotion skills using hierarchical deep reinforcement learning. *Acm transactions on graphics (tog)*, 36(4):1–13, 2017.
- Karl Pertsch, Youngwoon Lee, and Joseph Lim. Accelerating reinforcement learning with learned skill priors. In *Conference on robot learning*, pp. 188–204. PMLR, 2021.
- Doina Precup, Richard S Sutton, and Satinder Singh. Eligibility traces for off-policy policy evaluation. In *ICML*, volume 2000, pp. 759–766. Citeseer, 2000.

- Martin Riedmiller, Roland Hafner, Thomas Lampe, Michael Neunert, Jonas Degraeve, Tom Wiele, Vlad Mnih, Nicolas Heess, and Jost Tobias Springenberg. Learning by playing solving sparse reward tasks from scratch. In *International conference on machine learning*, pp. 4344–4353. PMLR, 2018.
- Mark Rowland, Will Dabney, and Rémi Munos. Adaptive trade-offs in off-policy learning. In *International Conference on Artificial Intelligence and Statistics*, pp. 34–44. PMLR, 2020.
- Younggyo Seo and Pieter Abbeel. Reinforcement learning with action sequence for data-efficient robot learning. 2024.
- Younggyo Seo, Jafar Uruç, and Stephen James. Continuous control with coarse-to-fine reinforcement learning. In *8th Annual Conference on Robot Learning*, 2024. URL <https://openreview.net/forum?id=WjDR48cL30>.
- Tanmay Shankar and Abhinav Gupta. Learning robot skills with temporal variational inference. In *International Conference on Machine Learning*, pp. 8624–8633. PMLR, 2020.
- Özgür Şimşek and Andrew G. Barto. Betweenness centrality as a basis for forming skills. Working-paper, University of Massachusetts Amherst, April 2007.
- Aravind Srinivas, Ramnandan Krishnamurthy, Peeyush Kumar, and Balaraman Ravindran. Option discovery in hierarchical reinforcement learning using spatio-temporal clustering. *arXiv preprint arXiv:1605.05359*, 2016.
- Richard S Sutton, Andrew G Barto, et al. *Reinforcement learning: An introduction*, volume 1. MIT press Cambridge, 1998.
- Richard S Sutton, Doina Precup, and Satinder Singh. Between MDPs and semi-MDPs: A framework for temporal abstraction in reinforcement learning. *Artificial intelligence*, 112(1-2):181–211, 1999.
- Denis Tarasov, Vladislav Kurenkov, Alexander Nikulin, and Sergey Kolesnikov. Revisiting the minimalist approach to offline reinforcement learning. *Advances in Neural Information Processing Systems*, 36, 2024.
- Philip S Thomas, Scott Niekum, Georgios Theodorou, and George Konidaris. Policy evaluation using the Ω -return. *Advances in Neural Information Processing Systems*, 28, 2015.
- Dong Tian, Ge Li, Hongyi Zhou, Onur Celik, and Gerhard Neumann. Chunking the critic: A transformer-based soft actor-critic with N-step returns. *arXiv preprint arXiv:2503.03660*, 2025.
- Ahmed Touati, Jérémy Rapin, and Yann Ollivier. Does zero-shot reinforcement learning exist? In *The Eleventh International Conference on Learning Representations*, 2022.
- Alexander Vezhnevets, Volodymyr Mnih, Simon Osindero, Alex Graves, Oriol Vinyals, John Agapiou, et al. Strategic attentive writer for learning macro-actions. *Advances in neural information processing systems*, 29, 2016.
- Alexander Sasha Vezhnevets, Simon Osindero, Tom Schaul, Nicolas Heess, Max Jaderberg, David Silver, and Koray Kavukcuoglu. Feudal networks for hierarchical reinforcement learning. In *International conference on machine learning*, pp. 3540–3549. PMLR, 2017.
- Max Wilcoxson, Qiyang Li, Kevin Frans, and Sergey Levine. Leveraging skills from unlabeled prior data for efficient online exploration. *arXiv preprint arXiv:2410.18076*, 2024.
- Kevin Xie, Homanga Bharadhwaj, Danijar Hafner, Animesh Garg, and Florian Shkurti. Latent skill planning for exploration and transfer. In *International Conference on Learning Representations*, 2021. URL <https://openreview.net/forum?id=jXe91kq3jAq>.
- Tony Z Zhao, Vikash Kumar, Sergey Levine, and Chelsea Finn. Learning fine-grained bimanual manipulation with low-cost hardware. *arXiv preprint arXiv:2304.13705*, 2023.

A FULL RESULTS

Table 2 reports the performance of our method (DQC) and baselines for all hyperparameter configurations. All of them use the same hyperparameters in Table 4 with the only exception that SHARSA handles goal-sampling for training behavior cloning policies separate as we discuss in more details in Appendix C.

Task	DQC ($h = 25$) ($h_a = 5$)	QC-NS ($n = 25$) ($h_a = 5$)	DQC ($h = 25$) ($h_a = 1$)	NS ($n = 25$)	QC ($h_a = 25$)	DQC ($h = 5$) ($h_a = 1$)	QC ($h_a = 5$)	NS ($n = 5$)	OS	SHARSA	HIQL	IQL	HFBC	FBC
cube-triple-100M	98 ^[96,98]	69 ^[13,98]	73 ^[67,80]	42 ^[5,74]	23 ^[0,50]	75 ^[70,79]	14 ^[6,28]	35 ^[0,81]	1 ^[0,2]	82 ^[78,88]	-	64 ^[59,68]	57 ^[54,61]	53 ^[48,57]
cube-quadruple-100M	90 ^[88,90]	37 ^[20,63]	46 ^[41,52]	20 ^[1,52]	2 ^[0,5]	58 ^[56,59]	24 ^[6,44]	27 ^[0,64]	0 ^[0,0]	67 ^[62,74]	-	53 ^[53,53]	38 ^[34,41]	32 ^[30,33]
cube-octuple-1B	24 ^[22,25]	28 ^[26,31]	13 ^[12,14]	8 ^[8,9]	1 ^[1,2]	0 ^[0,0]	0 ^[0,0]	0 ^[0,0]	0 ^[0,0]	20 ^[19,20]	1 ^[0,2]	0 ^[0,0]	20 ^[17,23]	0 ^[0,0]
humanoidmaze-giant	32 ^[29,36]	21 ^[19,24]	72 ^[67,75]	59 ^[54,64]	-	0 ^[0,0]	0 ^[0,0]	0 ^[0,0]	0 ^[0,1]	18 ^[13,25]	22 ^[18,29]	3 ^[3,4]	19 ^[16,22]	1 ^[0,2]
puzzle-4x5	96 ^[95,97]	97 ^[96,98]	89 ^[84,92]	91 ^[90,94]	-	19 ^[19,19]	20 ^[19,20]	19 ^[19,20]	14 ^[13,15]	1 ^[0,2]	5 ^[4,7]	19 ^[18,20]	4 ^[4,6]	0 ^[0,0]
puzzle-4x6-1B	65 ^[60,75]	69 ^[64,71]	72 ^[63,80]	74 ^[66,86]	-	23 ^[22,23]	20 ^[20,20]	22 ^[21,22]	17 ^[14,19]	56 ^[49,63]	8 ^[4,12]	17 ^[16,19]	2 ^[1,2]	0 ^[0,0]

Table 2: **Complete results for all configurations.** All means and 95% bootstrapped confidence intervals are computed over 3 seeds. (*) indicates that we take the results from the original paper (Park et al., 2025), where we take the results with larger 10M-sized datasets for humanoidmaze-giant (originally 4M) and puzzle-4x5 (originally 3M). We omit HIQL results on cube-{triple, quadruple} given its high computational cost and poor performance on the other tasks.

B ENVIRONMENTS AND DATASETS

To evaluate our method, we consider 8 goal-conditioned environments in OGBench with varying difficulties (Figure 4). The dataset size, episode length, and the action dimension for each environment is available in Table 3. We describe each of the environments and the datasets we use as follows.

Environment cube-*: We consider three cube environments (cube-triple, cube-quadruple, cube-octuple). As the names suggest, the goal of these environments involve using a robot arm to manipulate 3/4/8 cubes from some initial configuration to some specified goal configuration. We use the same five evaluation tasks used in OGBench (Park et al., 2024a) for cube-triple and cube-quadruple and the same five evaluation tasks used in Park et al. (2025) for cube-octuple. We refer the environment detail to the corresponding references.

Environment	Dataset Size	Episode Length	Action Dim. (A)
cube-triple-100M	100M	1000	5
cube-quadruple-100M	100M	1000	5
cube-octuple-1B	1B	1500	5
humanoidmaze-giant	4M (default)	4000	21
puzzle-4x5	3M (default)	1000	5
puzzle-4x6-1B	1B	1000	5

Table 3: **Environment metadata.** For both humanoidmaze-giant and puzzle-4x5, we use the default dataset that is released in the original OGBench benchmark (Park et al., 2024a). For the other environments, we use larger datasets as we find them to be essential for achieving good performances on these environments.

Environment humanoidmaze-*: We also consider the hardest locomotion environment available in OGBench. The goal of the environment is to control and navigate a humanoid agent from some initial location to some specified goal location in a 16×12 maze. This environment also has the longest episode length (4000, more than twice as long as the second longest episode length as used in cube-octuple). We refer the environment detail to Park et al. (2024a).

Environment puzzle-*: Finally, we consider two environments that involve solving a combinatorial puzzle with a robot arm. The puzzle consists of a board of 4×5 or 4×6 buttons, organized as a regular grid (4 rows and 5 or 6 columns). Each button has a binary state. Whenever the end-effector of the arm touches a button, the button and all its adjacent four buttons (three or two if the button is on the edge of the grid or in the corner) flip its binary state. The goal of the environment is to transform the board from some initial state to some specified goal state. We refer the environment detail to Park et al. (2025).

At the test-time/evaluation-time, the goal-conditioned agent is tested on five evaluation tasks for each of the six environments we consider. The overall success rate is the average over 5 tasks with 50 evaluation trials each.

Datasets. We use play datasets for all cube-* and puzzle-* environments and navigate dataset for humanoidmaze-*. We use the original datasets available for humanoidmaze-giant and puzzle-4x5 because they are sufficient for solving the environments. Using larger datasets on these environments do not help differentiating among different methods/baselines. For each of the other environments, we use the largest dataset available from [Park et al. \(2025\)](#) as we find it to be necessary to solve these environments (or achieve non-trivial performance on the hardest cube-octuple environment).

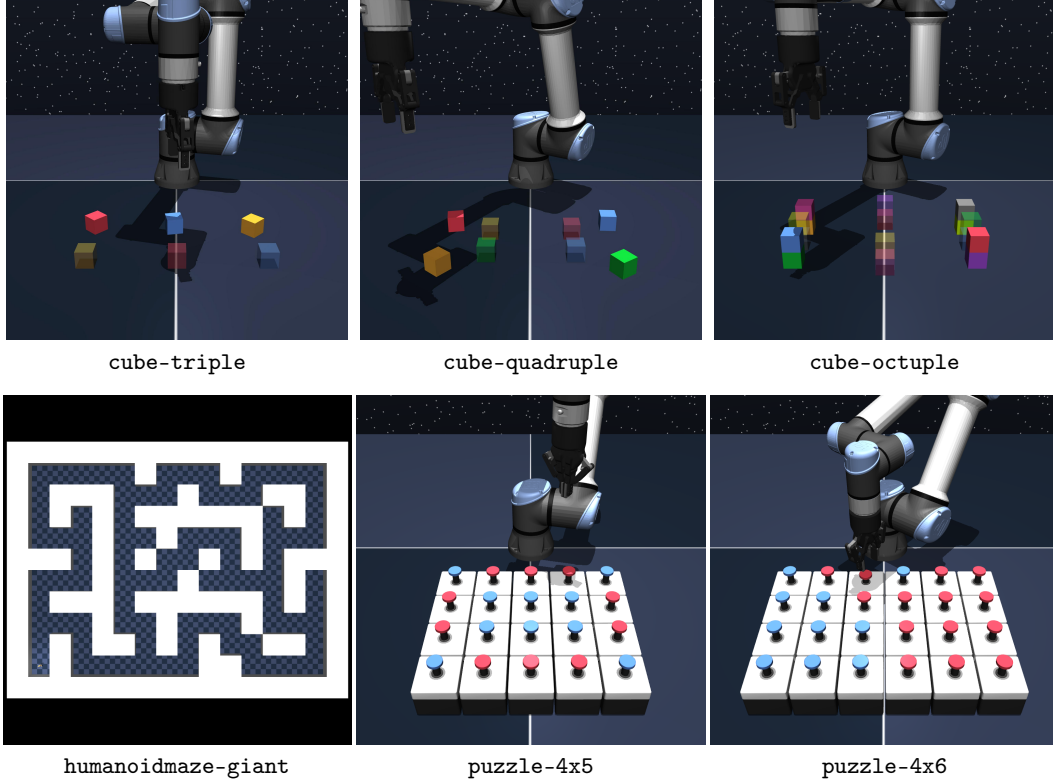


Figure 4: Visualization of environments.

C HYPERPARAMETERS AND IMPLEMENTATION DETAILS

Hyperparameters. Table 4 describes the common hyperparameters used in all our experiments (except for the ones with \dagger where the numbers are directly taken from prior work). Table 6 (for our method) and Table 7 (for baselines) describe the environment-specific hyperparameters.

Goal-conditioned RL implementation details. While we have described in the main body of the paper how DQC works as a general RL algorithm, we have not touched on how DQC and similarly all our baselines works with the goal-condition RL (GCRL) setting. We consider the setting where we have access to an oracle goal representation $\Psi : \mathcal{S} \rightarrow \mathcal{G}$ where \mathcal{G} is the goal space (see Table 5 for the oracle goal representation description for each environment). The goal-conditioned reward function $r : (s, g) \mapsto \mathbb{I}_{\Psi(s)=g}$ is a binary reward function where its output is 1 if the goal g is reached by the current state s . We can treat g as part of an extended state $\tilde{s} = [s, g] \in \tilde{\mathcal{S}} = \mathcal{S} \times \mathcal{G}$ and learn value functions (e.g., $Q_\phi(\tilde{s}, a)$) normally with such extended state.

A common trick in the GCRL setting is to use goal relabeling. That is, during training for each (s, a) pair in the training batch, a goal g is sampled from some distribution (i.e., $p^{\mathcal{D}}(\cdot \mid s, a)$) and the reward

Parameter	Value
Batch size	4096
Discount factor (γ)	0.999
Optimizer	Adam
Learning rate	3×10^{-4}
Target network update rate (λ)	5×10^{-3}
Critic ensemble size (K)	2
Critic target	$\min(Q_1, Q_2)$ for cube-* $(Q_1 + Q_2)/2$ for puzzle-* and humanoid-*
Implicit Backup Quantile (κ_b)	0.9
Value loss type	binary cross entropy
Best-of- N sampling (N)	32
Number of flow steps	10
Number of training steps	10^6
Network width	1024
Network depth	4 hidden layers
Value goal sampling ($w_{\text{cur}}^v, w_{\text{geom}}^v, w_{\text{traj}}^v, w_{\text{rand}}^v$)	(0.2, 0, 0.5, 0.3)
Actor goal sampling ($w_{\text{cur}}^p, w_{\text{geom}}^p, w_{\text{traj}}^p, w_{\text{rand}}^p$)	DQC/QC/NS/OS: π_β is not goal-conditioned SHARSA (cube): (0, 1, 0, 0) SHARSA (puzzle): (0, 0, 1, 0) SHARSA (humanoidmaze): (0, 0, 1, 0)

Table 4: **Common hyperparameters.** For the GCRL goal-sampling distribution we follow the same hyperparameters used in Park et al. (2025).

Environment	Goal Representation (Ψ)	Goal Domain (\mathcal{G})
cube-triple	(x, y, z) of three cubes (rel. to center)	\mathbb{R}^9
cube-quadruple	(x, y, z) of four cubes (rel. to center)	\mathbb{R}^{12}
cube-octuple	(x, y, z) of eight cubes (rel. to center)	\mathbb{R}^{24}
humanoidmaze-giant	(x, y) of the humanoid	\mathbb{R}^2
puzzle-4x5	the binary state for each button	$\{0, 1\}^{20}$
puzzle-4x6	the binary state for each button	$\{0, 1\}^{24}$

Table 5: **Oracle goal representation description for each environment.** Following Park et al. (2025), we assume access to an oracle goal representation for each environment. More detailed definition of these oracle goal representations is available in OGBench (Park et al., 2024a).

of the transition is relabeled with the goal-conditioned reward function. Following Park et al. (2025), the goal distribution $P^g(\cdot | s, a) : \mathcal{S} \times \mathcal{A} \rightarrow \Delta_{\mathcal{G}}$ is a mixture of four distributions, conditioned on the training state-action example:

$$P^g = w_{\text{cur}} P_{\text{cur}}^g + w_{\text{geom}} P_{\text{geom}}^g + w_{\text{traj}} P_{\text{traj}}^g + w_{\text{rand}} P_{\text{rand}}^g, \quad (28)$$

where

1. $P_{\text{cur}}^g(\cdot | s, a) = \delta_{\Psi(s)}$: the goal is the same as the current state;
2. $P_{\text{geom}}^g(\cdot | s, a)$: geometric distribution over the future states in the same trajectory that (s, a) is from;
3. $P_{\text{traj}}^g(\cdot | s, a)$: uniform distribution over the future states in the same trajectory that (s, a) is from; and finally
4. $P_{\text{rand}}^g(\cdot | s, a) = \Psi(\mathcal{U}_{\mathcal{D}(s)})$: uniform distribution over the dataset ($\mathcal{D}(s)$ is the distribution of states in the dataset).

and $w_{\text{cur}}, w_{\text{geom}}, w_{\text{traj}}, w_{\text{rand}} > 0$ are the corresponding weights for each of the distribution components with $w_{\text{cur}} + w_{\text{geom}} + w_{\text{traj}} + w_{\text{rand}} = 1$.

Environment	Distillation Expectile (κ_d)	Critic Chunk Size (h)	Policy Chunk Size (h_a)
cube-triple-100M	0.8	25	5
cube-quadruple-100M	0.8	25	5
cube-octuple-1B	0.8	25	5
humanoidmaze-giant	0.5	25	1
puzzle-4x5	0.5	25	5
puzzle-4x6-1B	0.5	25	1

Table 6: **Environment-specific hyperparameters for DQC.** We use a distillation expectile of 0.8 for all cube environments and a distillation expectile of 0.5 for all other environments. For all cube-* environments, we use a policy chunk size of 5 (*i.e.*, $h_a = 5$). For all other environments, we use 1-step policy (*i.e.*, $h_a = 1$). All environments use a critic chunk size of $h = 25$.

Environment	QC ($h = h_a$)	NS (n)	SHARSA (n)
cube-triple-100M	5	25	25
cube-quadruple-100M	5	25	25
cube-octuple-1B	5	25	25
humanoidmaze-giant-4M	5	25	50
puzzle-4x5-3M	5	25	50
puzzle-4x6-1B	5	25	50

Table 7: **Environment-specific hyperparameters for QC, NS, SHARSA .** For QC, we find $h = h_a = 5$ works the best for all environments. For NS, we find $n = 25$ works the best for all environments. For SHARSA, we follow the hyperparameters in the original paper (Park et al., 2025).

In practice, it has been found to be beneficial to use a separate set of goal sampling weights for TD backup (Park et al., 2024a) (*i.e.*, $(w_{\text{cur}}^v, w_{\text{geom}}^v, w_{\text{traj}}^v, w_{\text{rand}}^v)$) and for policy learning (*i.e.*, $(w_{\text{cur}}^p, w_{\text{geom}}^p, w_{\text{traj}}^p, w_{\text{rand}}^p)$). However, in our implementation of DQC/QC/NS/OS, we do not train a goal-conditioned policy, as our policy extraction is done entirely at test-time by best-of-N sampling from an *unconditional* (*i.e.*, not goal-conditioned) behavior policy π_β . In particular, we use an unconditioned flow policy $\pi_\beta(\cdot | s)$ that is parameterized by a velocity field $v_\beta : \mathcal{S} \times \mathbb{R}^A \times [0, 1] \rightarrow \mathbb{R}^A$ that is trained with the standard flow-matching objective:

$$L_{\text{FM}}(\beta) = \mathbb{E}_{u \sim \mathcal{U}[0,1], z \sim \mathcal{N}_{(s,a)} \sim \mathcal{D}} [\|v_\beta(s, (1-u)z + ua, u) - a + z\|_2^2] \quad (29)$$

For SHARSA, we use the official implementation where both flow policies (high-level and low-level) are goal-conditioned (and thus are trained with the goal distribution mixture specified by $w_{\text{cur}}^p, w_{\text{geom}}^p, w_{\text{traj}}^p, w_{\text{rand}}^p$). The goal sampling distribution for training the value networks (for all methods) and the goal sampling distribution for the policy networks (for SHARSA only) are provided in Table 4.

D EXAMPLES OF OPEN-LOOP CONSISTENT DATA

In this section, we provide some examples of open-loop consistent data that could serve as a useful basis for theoretical analyses in future work. The first example is any data collected from a near-deterministic dynamics, as formally defined as follows:

Definition D.1 (Near-deterministic Dynamics). *A transition dynamics T is ε -deterministic if there exists a deterministic transition dynamics represented by function $f : \mathcal{S} \times \mathcal{A} \rightarrow \mathcal{S}$ and another transition dynamics $\tilde{T} : \mathcal{S} \times \mathcal{A} \rightarrow \Delta_{\mathcal{S}}$, and T is a combination of f and \tilde{T} :*

$$T(s' | s, a) = (1 - \varepsilon)\delta_{f(s,a)}(s') + \varepsilon\tilde{T}(s' | s, a), \forall s, s' \in \mathcal{S}, a \in \mathcal{A}. \quad (30)$$

Theorem D.2 (Deterministic Dynamics are Open-loop Consistent). *If a transition dynamics \mathcal{M} is ε -deterministic, then any data \mathcal{D} is ε_h -open-loop consistent with respect to \mathcal{M} for any $h \in \mathbb{N}^+$ as long as $\varepsilon_h \geq 3(1 - (1 - \varepsilon)^{h-1})$.*

The proof of Theorem D.2 is available in Appendix E.

In addition to deterministic dynamics, any data collected by open-loop policies are also open-loop consistent.

Definition D.3 (Data Collected by Open-loop Policies). *A policy is open-loop if its action distribution does not depend on the state (i.e., $\pi(a_{t:t+h} \mid s_{t:t+h})$ is the same for all $s_{t:t+h}$).*

Remark D.4. *If the data \mathcal{D} is collected with an open-loop policy, then \mathcal{D} is strongly open-loop consistent.*

E PROOFS

Theorem 4.4 (Bias of Action Chunking Critic). *Let $\hat{V}_{ac} : \mathcal{S} \rightarrow [0, 1/(1 - \gamma)]$ be a solution of*

$$\hat{V}_{ac}(s_t) = \mathbb{E}_{s_{t+1:t+h+1}, a_{t:t+h} \sim P_{\mathcal{D}}(\cdot \mid s_t)} \left[R_{t:t+h} + \gamma^h \hat{V}_{ac}(s_{t+h}) \right], \quad (12)$$

with $R_{t:t+h} = \sum_{t'=t}^{t+h} \gamma^{t'-t} r(s_{t'}, a_{t'})$ and V_{ac} is the true value of $\pi_{\mathcal{D}}^{\text{open}} : s_t \mapsto P_{\mathcal{D}}(a_{t:t+h} \mid s_t)$. If \mathcal{D} is ε_h -open-loop consistent, then under $\text{supp}(\mathcal{D})$,

$$\|V_{ac} - \hat{V}_{ac}\|_{\infty} \leq \frac{\varepsilon_h}{(1 - \gamma^h)(1 - \gamma)}. \quad (13)$$

Proof. Since \mathcal{D} is $\varepsilon_{h'}$ -open-loop consistent in state-action for $h' < h$, the state-action distribution leading up to step h admits the following bound:

$$2D_{\text{TV}}(P_{\mathcal{D}}(s_{t+h}, a_{t+h} \mid s_t) \parallel P_{\mathcal{D}}^{\text{open}}(s_{t+h}, a_{t+h} \mid s_t)) \leq \varepsilon_h \quad (31)$$

Let $R_{t:t+h} = \sum_{k=0}^{h-1} \gamma^k r(s_{t+k}, a_{t+k})$ be the h -step reward distribution. Then the difference in h -step reward is bounded by

$$\left| \mathbb{E}_{P_{\mathcal{D}}(\cdot \mid s_t)}[R_{t:t+h}] - \mathbb{E}_{P_{\mathcal{D}}^{\text{open}}(\cdot \mid s_t)}[R_{t:t+h}] \right| \quad (32)$$

$$\leq \sum_{h'=0}^{h-1} \left[2\gamma^{h'} D_{\text{TV}}(P_{\mathcal{D}}(s_{t+h'}, a_{t+h'} \mid s_t) \parallel P_{\mathcal{D}}^{\text{open}}(s_{t+h'}, a_{t+h'} \mid s_t)) \right] \quad (33)$$

$$\leq \sum_{h'=0}^{h-1} \gamma^{h'} \varepsilon_h \quad (34)$$

Since \mathcal{D} is ε_h -open-loop consistent for h in state, we have

$$D_{\text{TV}}(P_{\mathcal{D}}(s_{t+h} \mid s_t) \parallel P_{\mathcal{D}}^{\text{open}}(s_{t+h} \mid s_t)) \leq \varepsilon_h \quad (35)$$

$$\left| \mathbb{E}_{s_{t+h} \sim P_{\mathcal{D}}(s_{t+h} \mid s_t)} [\hat{V}_{ac}(s_{t+h})] - \mathbb{E}_{s_{t+h} \sim P_{\mathcal{D}}^{\text{open}}(s_{t+h} \mid s_t)} [V_{ac}(s_{t+h})] \right| \quad (36)$$

$$\leq 2D_{\text{TV}}(P_{\mathcal{D}}(s_{t+h} \mid s_t) \parallel P_{\mathcal{D}}^{\text{open}}(s_{t+h} \mid s_t)) \mathbb{E}_{s_{t+h} \sim P_{\mathcal{D}}^{\text{open}}(\cdot \mid s_t)} [\hat{V}_{ac}(s_{t+h})] \quad (37)$$

$$+ \left| \mathbb{E}_{s_{t+h} \sim P_{\mathcal{D}}(s_{t+h} \mid s_t)} [\hat{V}_{ac}(s_{t+h}) - V_{ac}(s_{t+h})] \right| \quad (38)$$

$$\leq \frac{\varepsilon_h}{1 - \gamma} + \|\hat{V}_{ac} - V_{ac}\|_{\infty} \quad (39)$$

The first term in the last line works because even though $P_{\mathcal{D}}^{\text{open}}$ can go out of the support of \mathcal{D} , the value of \hat{V}_{ac} is uniformly bounded by its range $[0, 1/(1 - \gamma)]$. The second term works because the support of $s_{t+h} \mid s_t$ is a subset of the support for s_t as part of Assumption 4.1.

For all $s_t \in \text{supp}(P_{\mathcal{D}}(s_t))$,

$$\left| \hat{V}_{\text{ac}}(s_t) - V_{\text{ac}}(s_t) \right| \leq \left| \mathbb{E}_{P_{\mathcal{D}}(\cdot|s_t)}[R_{t:t+h}] - \mathbb{E}_{P_{\mathcal{D}}^{\text{open}}(\cdot|s_t)}[R_{t:t+h}] \right| \quad (40)$$

$$+ \gamma^h \left| \mathbb{E}_{s_{t+h} \sim P_{\mathcal{D}}(s_{t+h}|s_t)} [\hat{V}_{\text{ac}}(s_{t+h})] - \mathbb{E}_{s_{t+h} \sim P_{\mathcal{D}}^{\text{open}}(s_{t+h}|s_t)} [V_{\text{ac}}(s_{t+h})] \right| \quad (41)$$

$$\leq \sum_{h'=0}^{h-1} [\gamma^{h'} \varepsilon_h] + \frac{\gamma^h \varepsilon_h}{1-\gamma} + \gamma^h \|V_{\text{ac}} - \hat{V}_{\text{ac}}\|_{\infty}. \quad (42)$$

Therefore, under the same support,

$$\|V_{\text{ac}} - \hat{V}_{\text{ac}}\|_{\infty} \leq \frac{1}{1-\gamma^h} \left(\sum_{h'=0}^{h-1} [\gamma^{h'} \varepsilon_h] + \frac{\gamma^h \varepsilon_h}{1-\gamma} \right), \quad (43)$$

which can be simplified to be

$$\|V_{\text{ac}} - \hat{V}_{\text{ac}}\|_{\infty} \leq \frac{\varepsilon_h}{(1-\gamma)(1-\gamma^h)}. \quad (44)$$

□

Corollary 4.5 (Optimal Action Chunking Policy). *Let $\pi^* : \mathcal{S} \rightarrow \Delta_{\mathcal{A}}$ be an optimal policy in \mathcal{M} and \mathcal{D}^* be the data collected by π^* . If \mathcal{D}^* is ε_h -open-loop consistent, then under $\text{supp}(\mathcal{D}^*)$,*

$$\|V_{\text{ac}}^* - V^*\|_{\infty} \leq \|\tilde{V}_{\text{ac}} - V^*\|_{\infty} \leq \frac{\varepsilon_h}{(1-\gamma^h)(1-\gamma)}, \quad (14)$$

where V^* is the value of the optimal policy π^* , V_{ac}^* is the true value of the optimal action chunking policy, and \tilde{V}_{ac} is the true value of the action chunking policy from cloning the data \mathcal{D}^* :

$$\tilde{\pi}_{\text{ac}}(a_{t:t+h} | s_t) : s_t \mapsto P_{\mathcal{D}^*}(\cdot | s_t). \quad (15)$$

Proof. Let \hat{V}_{ac} be the fixed point of the following equation:

$$\hat{V}_{\text{ac}}(s_t) = \mathbb{E}_{s_{t+1:t+h+1}, a_{t:t+h} \sim P_{\mathcal{D}^*}(\cdot|s_t)} [R_{t:t+h} + \gamma^h \hat{V}_{\text{ac}}(s_{t+h})] \quad (45)$$

where again $R_{t:t+h} = \sum_{t'=t}^{t+h} \gamma^{t'-t} r(s_{t'}, a_{t'})$. The value of the optimal policy is the fixed point of the following equation:

$$V^*(s_t) = \mathbb{E}_{s_{t+1}, a_t \sim P_{\mathcal{D}^*}(\cdot|s_t)} [r(s_t, a_t) + \gamma V^*(s_{t+1})] \quad (46)$$

$$= \mathbb{E}_{s_{t:t+2}, a_{t:t+1} \sim P_{\mathcal{D}^*}(\cdot|s_t)} [r(s_t, a_t) + \gamma r(s_{t+1}, a_{t+1}) + \gamma V^*(s_{t+2})] \quad (47)$$

$$\dots \quad (48)$$

$$= \mathbb{E}_{s_{t+1:t+h+1}, a_{t:t+h} \sim P_{\mathcal{D}^*}(\cdot|s_t)} [R_{t:t+h} + \gamma^h V^*(s_{t+h})] \quad (49)$$

which is equivalent to fixed-point equation for \hat{V}_{ac} . Therefore $\hat{V}_{\text{ac}} = V^*$. By Theorem 4.4, we know that the true value V_{ac} of the action chunking policy $\tilde{\pi}_{\text{ac}}$ that clones \mathcal{D}^* is close to \hat{V}_{ac} :

$$\|\hat{V}_{\text{ac}} - \tilde{V}_{\text{ac}}\|_{\infty} \leq \frac{\varepsilon_h}{(1-\gamma^h)(1-\gamma)} \quad (50)$$

which means that

$$\|V^* - \tilde{V}_{\text{ac}}\|_{\infty} \leq \frac{\varepsilon_h}{(1-\gamma^h)(1-\gamma)} \quad (51)$$

Since the optimal action chunking policy, by definition, attains equally good or better values (over \mathcal{S}) represented by V_{ac} , and the optimal policy π^* also attains equally good or better value (i.e., V^*) compared to that of the optimal action chunking policy π_{ac}^* (i.e., V_{ac}^*), the following inequality holds under $\text{supp}(\mathcal{D}^*)$:

$$V^* \geq V_{\text{ac}}^* \geq \tilde{V}_{\text{ac}}. \quad (52)$$

Therefore,

$$\|V_{\text{ac}}^* - V^*\| \leq \|\tilde{V}_{\text{ac}} - V^*\| \leq \frac{\varepsilon_h}{(1-\gamma^h)(1-\gamma)}$$

□

Theorem 4.6 (Q-Learning with Action Chunking Policy on Off-policy Data). *If \mathcal{D} is strongly ε_h -open-loop consistent and $\text{supp}(\mathcal{D}) \supseteq \text{supp}(\mathcal{D}^*)$, with \mathcal{D}^* being the data distribution of an arbitrary optimal policy π^* under \mathcal{M} , then the following bound holds under $\text{supp}(\mathcal{D}^*)$:*

$$\|V_{\text{ac}}^+ - V^*\|_\infty \leq \frac{3\varepsilon_h}{(1-\gamma^h)(1-\gamma)}, \quad (18)$$

where V^* is the value of an optimal policy under \mathcal{M} .

Proof. We start by constructing a bound between \hat{Q}_{ac}^+ and Q_{ac}^* , the solution of the following bellman equation:

$$Q_{\text{ac}}^*(s_t, a_{t:t+h}) = \mathbb{E}_{s_{t+1:t+h+1} \sim P_{\mathcal{D}}^{\text{open}}(\cdot | s_t, a_{t:t+h})} \left[R_{t:t+h} + \gamma^h \max_{a_{t+h:t+2h}} Q_{\text{ac}}^*(s_{t+h}, a_{t+h:t+2h}) \right]. \quad (53)$$

Intuitively, Q_{ac}^* is the Q-function of the optimal action chunking policy π_{ac}^* that can be learned from \mathcal{D} . Because $\text{supp}(\mathcal{D}) \supseteq \text{supp}(\mathcal{D}^*)$, π_{ac}^* is at least as good as $\tilde{\pi}_{\text{ac}}$, the action chunking policy obtained by behavior cloning \mathcal{D}^* . Bounding the difference between \hat{Q}_{ac}^+ and Q_{ac}^* allows us to leverage the bound in Corollary 4.5 to form a bound between \hat{V}_{ac}^+ and V^* .

Since \mathcal{D} is strongly ε_h -open-loop consistent,

$$D_{\text{TV}}(P_{\mathcal{D}}^{\text{open}}(s_{t+h'} | s_t, a_{t:t+h'}) \| P_{\mathcal{D}}(s_{t+h'} | s_t, a_{t:t+h'})), \forall h' \in \{0, 1, \dots, h-1\} \quad (54)$$

Now, for the h -step reward, we have

$$\left| \mathbb{E}_{P_{\mathcal{D}}(\cdot | s_t, a_{t:t+h})} [R_{t:t+h}] - \mathbb{E}_{P_{\mathcal{D}}^{\text{open}}(\cdot | s_t, a_{t:t+h})} [R_{t:t+h}] \right| \quad (55)$$

$$\leq \sum_{h'=0}^{h-1} [D_{\text{TV}}(P_{\mathcal{D}}^{\text{open}}(s_{t+h'} | s_t, a_{t:t+h'}) \| P_{\mathcal{D}}(s_{t+h'} | s_t, a_{t:t+h'}))] \quad (56)$$

$$\leq \frac{(1-\gamma^h)\varepsilon_h}{1-\gamma} \quad (57)$$

Similarly, for the value h -step into the future, we have

$$\left| \mathbb{E}_{s_{t+h} \sim P_{\mathcal{D}}^{\text{open}}(s_{t+h} | s_t)} [V_{\text{ac}}^*(s_{t+h}, a_{t+h:t+2h})] - \mathbb{E}_{s_{t+h} \sim P_{\mathcal{D}}(s_{t+h} | s_t)} [\hat{V}_{\text{ac}}^+(s_{t+h}, a_{t+h:t+2h})] \right| \quad (58)$$

$$\leq 2D_{\text{TV}}(P_{\mathcal{D}}^{\text{open}}(s_{t+h'} | s_t, a_{t:t+h'}) \| P_{\mathcal{D}}(s_{t+h'} | s_t, a_{t:t+h'})) \|V_{\text{ac}}^*\|_\infty + \|V_{\text{ac}}^* - \hat{V}_{\text{ac}}^+\|_\infty \quad (59)$$

$$\leq \frac{\varepsilon_h}{1-\gamma} + \|V_{\text{ac}}^* - \hat{V}_{\text{ac}}^+\|_\infty \quad (60)$$

Combining the bound for the h -step reward and the bound on the value for s_{t+h} , we get

$$|Q_{\text{ac}}^*(s_t, a_{t:t+h}) - \hat{Q}_{\text{ac}}^+(s_t, a_{t:t+h})| \leq \frac{\varepsilon_h}{1-\gamma} + \gamma^h \|V_{\text{ac}}^* - \hat{V}_{\text{ac}}^+\|_\infty, \quad (61)$$

which can be recursively expanded to get

$$\|V_{\text{ac}}^* - \hat{V}_{\text{ac}}^+\|_\infty \leq \frac{\varepsilon_h}{(1-\gamma)(1-\gamma^h)} \quad (62)$$

By Theorem 4.4, we have

$$\|\hat{V}_{\text{ac}}^+ - V_{\text{ac}}^+\|_\infty \leq \frac{\varepsilon_h}{(1-\gamma)(1-\gamma^h)} \quad (63)$$

By Corollary 4.5, we have

$$\|V^* - \tilde{V}_{\text{ac}}\|_\infty \leq \frac{\varepsilon_h}{(1-\gamma)(1-\gamma^h)} \quad (64)$$

Since $\text{supp}(\mathcal{D}) \supseteq \text{supp}(\mathcal{D}^*)$, we know that V_{ac}^* is at least as good as \tilde{V}_{ac} uniformly.

Combining the three inequalities above, we get

$$\|V^* - V_{\text{ac}}^+\|_\infty \leq \frac{3\varepsilon_h}{(1-\gamma)(1-\gamma^h)} \quad (65)$$

□

Theorem 4.8. *Let \mathcal{D} be strongly ε_h -open-consistent, δ_n -suboptimal, and $\text{supp}(\mathcal{D}) \supseteq \text{supp}(\mathcal{D}^*)$. Let π_n^* be the optimal n -step return policy learned from \mathcal{D} , as the solution of*

$$Q_n^*(s_t, a_t) = \mathbb{E}_{P_{\mathcal{D}}} [R_{t:t+n} + \gamma^n Q_n^*(s_{t+n}, \pi_n^*(s_{t+n}))], \quad \pi_n^* : s_t \mapsto \arg \max_{a_t} Q_n^*(s_t, a_t). \quad (20)$$

As long as $\delta_n > \frac{3\varepsilon_h(1-\gamma^n)}{(1-\gamma)(1-\gamma^h)}$, then from all $s \in \text{supp}(\mathcal{D}^)$, the action chunking policy, π_{ac}^+ (Equation (17)), is better than the n -step return policy, π_n (Equation (20)) (i.e., $V_{\text{ac}}^+(s) > V_n^*(s)$).*

To prove Theorem 4.8, we first prove the following helper Lemma E.1 to quantify sub-optimality for n -step return policy.

Lemma E.1. *Let Q_n^* be the solution of the uncorrected n -step return backup equation:*

$$Q_n^*(s_t, a_t) = \mathbb{E}_{P_{\mathcal{D}}(\cdot|s_t, a_t)} \left[R_{t:t+n} + \gamma^n \max_{a_{t+n}} Q_n^*(s_{t+n}, a_{t+n}) \right] \quad (66)$$

The following inequality holds as long as \mathcal{D} is δ_n -suboptimal:

$$Q^*(s_t, a_t) \geq Q_n^*(s_t, a_t) + \frac{\delta_n}{1-\gamma^n}, \forall s_t \in \mathcal{S}, a_t \in \mathcal{A} \quad (67)$$

where Q^ is the Q -function of the optimal policy in \mathcal{M} . For the n -step return policy*

$$\pi_n^* : s_t \mapsto \arg \max_{a_t} Q_n^*(s_t, a_t), \quad (68)$$

its corresponding value admits a similar bound:

$$V^*(s_t) \geq V_n^*(s_t) + \frac{\delta_n}{1-\gamma^n}, \forall s_t \quad (69)$$

Proof. Using the definition of suboptimal data (Definition 4.7), we have

$$Q_n^*(s_t, a_t) = \mathbb{E}_{P_{\mathcal{D}}(\cdot|s_t, a_t)} \left[R_{t:t+n} + \gamma^n \max_{a_{t+n}} Q_n^*(s_{t+n}, a_{t+n}) \right] \quad (70)$$

$$\leq Q^*(s_t, a_t) - \delta_n + \gamma^n \mathbb{E}_{P_{\mathcal{D}}(\cdot|s_t, a_t)} \left[\max_{a_{t+n}} Q_n^*(s_{t+n}, a_{t+n}) - V^*(s_{t+h}) \right] \quad (71)$$

Rearranging the inequality above yields

$$Q_n^*(s_t, a_t) - Q^*(s_t, a_t) \leq -\delta_n + \gamma^n \mathbb{E}_{P_{\mathcal{D}}(\cdot|s_t)} [V_n^*(s_{t+n}) - V^*(s_{t+n})], \forall s_t \in \mathcal{S}, a_t \in \mathcal{A} \quad (72)$$

By recursively applying the inequality above, we have

$$Q^*(s_t, a_t) \geq Q_n^*(s_t, a_t) + \frac{\delta_n}{1-\gamma^n}, \forall s_t \in \mathcal{S}, a_t \in \mathcal{A} \quad (73)$$

By choosing $a_t^* = \pi_n^*(s_t)$, we see that

$$V^*(s_t) \geq Q^*(s_t, a_t) \quad (74)$$

$$\geq Q_n^*(s_t, a_t^*) + \frac{\delta_n}{1-\gamma^n} \quad (75)$$

$$= V_n^*(s_t) + \frac{\delta_n}{1-\gamma^n} \quad (76)$$

□

Now we are ready to prove the main Theorem 4.8.

Proof of Theorem 4.8. From Lemma E.1 and Theorem 4.6, we have

$$V_n^*(s) + \frac{\delta_n}{1 - \gamma^n} \leq V^*(s) \leq V_{ac}^+(s) + \frac{3\varepsilon_h}{(1 - \gamma^h)(1 - \gamma)} \quad (77)$$

Rearranging the terms give

$$V_{ac}^+(s) - V_n^*(s) \geq \frac{\delta_n}{1 - \gamma^n} - \frac{3\varepsilon_h}{(1 - \gamma)(1 - \gamma^h)} > 0 \quad (78)$$

□

Theorem D.2 (Deterministic Dynamics are Open-loop Consistent). *If a transition dynamics \mathcal{M} is ε -deterministic, then any data \mathcal{D} is ε_h -open-loop consistent with respect to \mathcal{M} for any $h \in \mathbb{N}^+$ as long as $\varepsilon_h \geq 3(1 - (1 - \varepsilon)^{h-1})$.*

Proof. Since T is ε -deterministic, it can be represented as $T(\cdot | s, a) = (1 - \varepsilon)\delta_{f(s,a)} + \varepsilon\tilde{T}(\cdot | s, a)$ for some $f : \mathcal{S} \times \mathcal{A} \rightarrow \mathcal{S}$ and $\tilde{T} : \mathcal{S} \times \mathcal{A} \rightarrow \Delta_{\mathcal{S}}$. Let $f(s, a_1, \dots, a_h) = f(\dots f(f(s, a_1), a_2) \dots a_h)$

Let $I \in \{0, 1\}$ a binary indicator variable that is 1 if and only if

$$s_{t+k+1} = f(s_{t+k}, a_{t+k}), \forall k \in \{0, 1, 2, \dots, h-1\} \quad (79)$$

Intuitively $I = 1$ when the trajectory is generated deterministically until but not including the last state s_h in the trajectory chunk.

From the fact that T is ε -deterministic, we know that

$$P_{\mathcal{D}}(I_h = 1) \geq (1 - \varepsilon)^{h-1} \quad (80)$$

We also have

$$P_{\mathcal{D}}(a_{t:t+h} | s_t) = P_{\mathcal{D}}(I_h = 1)P_{\mathcal{D}}(a_{t:t+h} | s_t, I_h = 1) + P_{\mathcal{D}}(I_h = 0)P_{\mathcal{D}}(a_{t:t+h} | s_t, I_h = 0) \quad (81)$$

Then we have

$$2D_{TV}(P_{\mathcal{D}}(a_{t:t+h} | s_t) || P_{\mathcal{D}}(a_{t:t+h} | s_t, I_h = 1)) \leq (1 - (1 - \varepsilon)^{h-1}) \quad (82)$$

If we transform each distribution of $a_{t:t+h}$ deterministically by $f(s_t, \cdot)$, by data processing inequality we have

$$2D_{TV}(\mathbb{E}_{a_{t:t+h} \sim P_{\mathcal{D}}(\cdot | s_t)} [\delta_{f(s_t, a_{t:t+h})}] || \mathbb{E}_{a_{t:t+h} \sim P_{\mathcal{D}}(\cdot | s_t, I_h = 1)} [\delta_{f(s_t, a_{t:t+h})}]) \leq (1 - (1 - \varepsilon)^{h-1}) \quad (83)$$

Similarly, we have

$$2D_{TV}(P_{\mathcal{D}}(a_{t:t+h+1} | s_t) || P_{\mathcal{D}}(a_{t:t+h+1} | s_t, I_{h+1} = 1)) \leq (1 - (1 - \varepsilon)^h) \quad (84)$$

which can be also deterministically transformed by taking $a_{t:t+h+1} \mapsto (f(s_t, \cdot), a_{t+h})$ to obtain

$$2D_{TV}(\mathbb{E}_{a_{t:t+h} \sim P_{\mathcal{D}}(\cdot | s_t)} [\pi_{\mathcal{D}}^{\text{open}}(a_{t+h} | s_t, a_{t:t+h}) \mathbb{I}_{f(s_t, a_{t:t+h})}] || \quad (85)$$

$$\mathbb{E}_{a_{t:t+h} \sim P_{\mathcal{D}}(\cdot | s_t, I_{h+1} = 1)} [\pi_{\mathcal{D}}^{\text{open}}(a_{t+h} | s_t, a_{t:t+h}, I_{h+1} = 1) \mathbb{I}_{f(s_t, a_{t:t+h})}]) \leq (1 - (1 - \varepsilon)^h) \quad (86)$$

Now, if we analyze the distribution of s_{t+h} subject to the open-loop execution of the action sequence from $P_{\mathcal{D}}(\cdot | s_t)$ and break it up into the deterministic and the non-deterministic case, we get

$$\mathbb{E}_{a_{t:t+h} \sim P_{\mathcal{D}}(\cdot | s_t)} [T_{a_{t:t+h}}(\cdot | s_t)] = P_T(I = 1)\mathbb{E}_{a_{t:t+h} \sim P_{\mathcal{D}}(\cdot | s_t)} [\delta_{f(s_t, a_{t:t+h})}] + \quad (87)$$

$$P_T(I = 0)\mathbb{E}_{a_{t:t+h} \sim P_{\mathcal{D}}(\cdot | s_t)} [T_{a_{t:t+h}}(\cdot | s_t, I_h = 0)] \quad (88)$$

Note that $P_T(I = 1)$ denotes the probability that an open-loop executed trajectory using $a_{t:t+h} \sim P_{\mathcal{D}}(\cdot | s_t)$ is deterministic. This is different from $P_{\mathcal{D}}(I_h = 1)$ because the latter is based on

$P_{\mathcal{D}}(s_{t:t+h+1}, a_{t:t+h})$ whereas $P_T(I_h = 1)$ is based on the open-loop trajectory distribution: $P_{\mathcal{D}}(\cdot | s_t) \prod_{k=0}^{h-1} T(s_{t+k} | s_t, a_{t:t+k})$. They both admit the same lower bound of $(1 - (1 - \varepsilon)^{h-1})$.

Therefore,

$$2D_{\text{TV}}(\mathbb{E}_{a_{t:t+h} \sim P_{\mathcal{D}}(\cdot | s_t)} [T_{a_{t:t+h}}(\cdot | s_t)] \parallel \mathbb{E}_{a_{t:t+h} \sim P_{\mathcal{D}}(\cdot | s_t)} [\delta_{f(s_t, a_{t:t+h})}]) \leq (1 - (1 - \varepsilon)^{h-1}) \quad (89)$$

Similarly for the state-action case, we can multiply both side by the same conditional distribution $\pi_{\mathcal{D}}^{\text{open}}(a_{t+h} | s_t, a_{t:t+h})$ which preserves the TV bound. For the left-hand side, we have

$$P_{\mathcal{D}}^{\text{open}}(s_{t+h}, a_{t+h} | s_t) = \mathbb{E}_{a_{t:t+h} \sim P_{\mathcal{D}}(\cdot | s_t)} [\pi_{\mathcal{D}}^{\text{open}}(a_{t+h} | s_t, a_{t:t+h}) T_{a_{t:t+h}}(s_{t+h} | s_t)] \quad (90)$$

Therefore, we get

$$2D_{\text{TV}}(P_{\mathcal{D}}^{\text{open}}(s_{t+h}, a_{t+h} | s_t) \parallel \mathbb{E}_{a_{t:t+h} \sim P_{\mathcal{D}}(\cdot | s_t)} [\pi_{\mathcal{D}}^{\text{open}}(a_{t+h} | s_t, a_{t:t+h}) \mathbb{I}_{f(s_t, a_{t:t+h})}]) \quad (91)$$

$$\leq (1 - (1 - \varepsilon)^{h-1}) \quad (92)$$

We also have

$$P_{\mathcal{D}}(s_{t+h} | s_t) = (1 - \varepsilon)^{h-1} P_{\mathcal{D}}(s_{t+h} | s_t, I = 1) + (1 - (1 - \varepsilon)^{h-1}) P_{\mathcal{D}}(s_{t+h} | s_t, I_h = 0) \quad (93)$$

Similarly, we have

$$2D_{\text{TV}}(P_{\mathcal{D}}(s_{t+h} | s_t) \parallel P_{\mathcal{D}}(s_{t+h} | s_t, I_h = 1)) \quad (94)$$

$$= 2D_{\text{TV}}(P_{\mathcal{D}}(s_{t+h} | s_t) \parallel \mathbb{E}_{a_{t:t+h} \sim P_{\mathcal{D}}(\cdot | s_t, I_h = 1)} [\delta_{f(s_t, a_{t:t+h})}]) \leq (1 - (1 - \varepsilon)^{h-1}) \quad (95)$$

For state-action, we can also get

$$P_{\mathcal{D}}(s_{t+h}, a_{t+h} | s_t) = (1 - \varepsilon)^h P_{\mathcal{D}}(s_{t+h}, a_{t+h} | s_t, I_{h+1} = 1) \quad (96)$$

$$+ (1 - (1 - \varepsilon)^h) P_{\mathcal{D}}(s_{t+h}, a_{t+h} | s_t, I_{h+1} = 0) \quad (97)$$

which can be turned into the TV distance bound:

$$2D_{\text{TV}}(P_{\mathcal{D}}(s_{t+h}, a_{t+h} | s_t) \parallel P_{\mathcal{D}}(s_{t+h}, a_{t+h} | s_t, I_{h+1} = 1)) \quad (98)$$

$$= 2D_{\text{TV}}(P_{\mathcal{D}}(s_{t+h}, a_{t+h} | s_t) \parallel \mathbb{E}_{a_{t:t+h} \sim P_{\mathcal{D}}(\cdot | s_t, I_{h+1} = 1)} [\pi_{\mathcal{D}}^{\text{open}}(a_{t+h} | s_t, a_{t:t+h}, I_{h+1} = 1) \mathbb{I}_{f(s_t, a_{t:t+h})}]) \quad (99)$$

$$\leq (1 - (1 - \varepsilon)^h) \quad (100)$$

$$\leq (1 - (1 - \varepsilon)^h) \quad (101)$$

Connecting all three total variation inequality (Equations (83), (89) and (94)) together, we get

$$2D_{\text{TV}}(P_{\mathcal{D}}(s_{t+h} | s_t) \parallel \mathbb{E}_{a_{t:t+h} \sim P_{\mathcal{D}}(\cdot | s_t)} [T_{a_{t:t+h}}(\cdot | s_t)]) \leq 3(1 - (1 - \varepsilon)^{h-1}) \leq \varepsilon_h \quad (102)$$

Connecting all three total variable inequality for state-action (Equations (85), (90) and (98)) together, we get

$$2D_{\text{TV}}(P_{\mathcal{D}}^{\text{open}}(s_{t+h-1}, a_{t+h-1} | s_t) \parallel P_{\mathcal{D}}(s_{t+h}, a_{t+h} | s_t)) \leq 3 - 2(1 - \varepsilon)^{h-1} - (1 - \varepsilon)^{h-2} \quad (103)$$

$$\leq 3(1 - (1 - \varepsilon)^{h-1}) \quad (104)$$

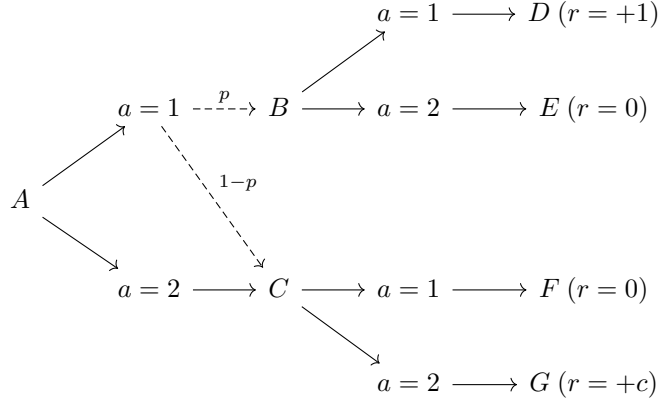
$$\leq \varepsilon_h \quad (105)$$

Therefore, \mathcal{D} is ε_h -open-loop consistent as desired. \square

F A PATHOLOGICAL FAILURE OF ACTION CHUNKING POLICIES WITHOUT THE STRONG OPEN-LOOP CONSISTENCY ASSUMPTION

In this section, we show an example where the optimal action chunking policy defined in Equation (17) can be highly suboptimal in the absence of the strong open-loop consistency condition.

We define an MDP as follows. Let $\mathcal{S} = \{A, B, C, D, E, F, G\}$ and $\mathcal{A} = \{1, 2\}$. Define the transition dynamics and reward function as shown in the diagram below:



where $p, c \in (0, 1)$ are real numbers and dotted lines denote stochastic transitions. For simplicity, assume that the MDP has a length-2 finite horizon with $\gamma = 1$, and the reward function depends only on states ($r(A) = r(B) = r(C) = r(E) = r(F) = 0$, $r(D) = 1$, and $r(G) = c$). Assume that the dataset is collected by a policy $\pi_{\mathcal{D}}$ defined as $\pi_{\mathcal{D}}(A) = 1$ (with probability 0.5) or 2 (with probability 0.5), $\pi_{\mathcal{D}}(B) = 1$ (with probability 1), and $\pi_{\mathcal{D}}(C) = 2$ (with probability 1).

Then, we have the following:

$$P_{\mathcal{D}}(A, (1, 1)) = D, \quad R(A, (1, 1)) = 1, \quad (106)$$

$$P_{\mathcal{D}}(A, (1, 2)) = G, \quad R(A, (1, 2)) = c, \quad (107)$$

$$P_{\mathcal{D}}(A, (2, 2)) = G, \quad R(A, (2, 2)) = c, \quad (108)$$

where we denote action chunks as a tuple and slightly abuse notation to denote deterministic outputs of $P_{\mathcal{D}}(\cdot \mid s_0, a_{0:2})$ (e.g., $P_{\mathcal{D}}(A, (1, 1)) = D$ indicates that all length-2 trajectories in \mathcal{D} from state A with $a_0 = a_1 = 1$ have $s_2 = D$ with probability 1). From this, we can compute \hat{Q}_{ac}^+ as follows:

$$\hat{Q}_{\text{ac}}^+(A, (1, 1)) = 1, \quad (109)$$

$$\hat{Q}_{\text{ac}}^+(A, (1, 2)) = c, \quad (110)$$

$$\hat{Q}_{\text{ac}}^+(A, (2, 2)) = c. \quad (111)$$

Then, assuming the missing data has a Q-value of 0 (i.e., $\hat{Q}_{\text{ac}}^+(A, (2, 1)) = 0$), the optimal action chunking policy is defined as $\hat{\pi}_{\text{ac}}^+(A) = (1, 1)$ (Equation (17)).

The true value of this action chunking policy is p . However, if p is small enough and c is large enough, the optimal strategy in this MDP is to always choose $(a_0, a_1) = (2, 2)$, in which case the agent receives a constant return of c . The suboptimality in this example is therefore $c - p$, which can be made arbitrarily close to 1 (the maximum possible regret in any finite, length-2 sparse-reward MDP with a terminal reward bounded by $[0, 1]$). This shows a pathological failure of an action chunking policy without the strong open-loop consistency assumption.

G ADDITIONAL RELATED WORK ON HIERARCHICAL REINFORCEMENT LEARNING

Hierarchical reinforcement learning methods (Dayan & Hinton, 1992; Dietterich, 2000; Peng et al., 2017; Riedmiller et al., 2018; Shankar & Gupta, 2020; Pertsch et al., 2021; Gehring et al., 2021; Xie et al., 2021) solve tasks by typically leveraging a bi-level structure: a set of low-level/skill policies that directly interact with the environment and a high-level policy that selects among low-level policies. The low-level policies can also be learned via online RL (Kulkarni et al., 2016; Vezhnevets et al., 2016; 2017; Nachum et al., 2018) or offline pre-training on a prior dataset (Paraschos et al., 2013;

Merel et al., 2018; Ajay et al., 2021; Pertsch et al., 2021; Touati et al., 2022; Nasiriany et al., 2022; Hu et al., 2023; Frans et al., 2024; Chen et al., 2024; Park et al., 2024b). In the options framework, these low-level policies are often additionally associated with initiation and termination conditions that specify when and for how long these actions can be used (Sutton et al., 1999; Menache et al., 2002; Chentanez et al., 2004; Şimşek & Barto, 2007; Konidaris, 2011; Daniel et al., 2016; Srinivas et al., 2016; Fox et al., 2017; Bacon et al., 2017; Bagaria & Konidaris, 2019; Bagaria et al., 2024; de Mello Koch et al., 2025). A long-lasting challenge in HRL is optimization stability because the high-level policy needs to optimize for an objective that is shaped by the constantly changing low-level policies (Nachum et al., 2018). Prior work (Ajay et al., 2021; Pertsch et al., 2021; Wilcoxson et al., 2024) avoided this by first pre-train low-level policies and then keep them frozen during the optimization of the high-level policy. Macro-actions (McGovern & Sutton, 1998; Durugkar et al., 2016), or action chunking (Zhao et al., 2023) is another form of temporally extended action, a special case of the low-level policies often considered in HRL, options literature, where a short horizon of actions are predicted all at once and executed in open loop. Such approach collapses the bi-level structure, conveniently side stepping optimization instability, and when combined with Q-learning, has shown great empirical successes in offline-to-online RL setting (Seo et al., 2024; Li et al., 2025b). Action chunking policies need to predict multiple actions open-loop, which can be difficult to learn and sacrifice reactivity. Our approach regains policy reactivity by predicting and executing only a partial action chunk, while still learning with the fully chunked critic for TD-backup. This design preserves the value propagation benefits of chunked critic without relying on fully open-loop action chunking policies, allowing our approach to work well on a wider range of tasks.