

# An Investigation of Offline Reinforcement Learning in Factorisable Action Spaces

Alex Beeson<sup>1,2</sup>, David Ireland<sup>1</sup>, Giovanni Montana<sup>1,3,4</sup>

<sup>1</sup>Warwick Manufacturing Group, University of Warwick, Coventry, UK

<sup>2</sup>Warwick Medical School, University of Warwick, Coventry, UK

<sup>3</sup>Department of Statistics, University of Warwick, Coventry, UK

<sup>4</sup>Alan Turing Institute, London, UK

*alex.beeson@warwick.ac.uk*

*david.ireland@warwick.ac.uk*

*g.montana@warwick.ac.uk*

Reviewed on OpenReview: <https://openreview.net/forum?id=STuxyUfpNV>

## Abstract

Expanding reinforcement learning (RL) to offline domains generates promising prospects, particularly in sectors where data collection poses substantial challenges or risks. Pivotal to the success of transferring RL offline is mitigating overestimation bias in value estimates for state-action pairs absent from data. Whilst numerous approaches have been proposed in recent years, these tend to focus primarily on continuous or small-scale discrete action spaces. Factorised discrete action spaces, on the other hand, have received relatively little attention, despite many real-world problems naturally having factorisable actions. In this work, we undertake a formative investigation into offline reinforcement learning in factorisable action spaces. Using value-decomposition as formulated in DecQN as a foundation, we present the case for a factorised approach and conduct an extensive empirical evaluation of several offline techniques adapted to the factorised setting. In the absence of established benchmarks, we introduce a suite of our own comprising datasets of varying quality and task complexity. Advocating for reproducible research and innovation, we make all datasets available for public use alongside our code base.

## 1 Introduction

The idea of transferring the successes of reinforcement learning (RL) to the offline setting is an enticing one. The opportunity for agents to learn optimal behaviour from sub-optimal data prior to environment interaction extends RL’s applicability to domains where data collection is costly, time-consuming or dangerous (Lange et al., 2012). This includes not only those domains where RL has traditionally found favour, such as games and robotics (Mnih et al., 2013; Hessel et al., 2018; Kalashnikov et al., 2018; Mahmood et al., 2018), but also areas in which online learning presents significant practical and/or ethical challenges, such as autonomous driving (Kiran et al., 2022) and healthcare (Yu et al., 2021a).

Unfortunately, taking RL offline is not as simple as naively applying standard off-policy algorithms to pre-existing datasets. A substantial challenge arises from the compounding and propagation of overestimation bias in value estimates for actions absent from data (Fujimoto et al., 2019b). This bias stems from the underlying bootstrapping procedure used to derive such estimates and subsequent maximisation to obtain policies, whether implicit such as in Q-learning or explicit as per actor-critic methods (Levine et al., 2020). Fundamentally, agents find it difficult to accurately determine the value of actions not previously encountered, and thus any attempt to determine optimal behaviour based on these values is destined to fail. Online, such

inaccuracies can be compensated for through continual assimilation of environmental feedback, but offline such a corrective mechanism is no longer available.

In response to these challenges, there has been a plethora of approaches put forward that aim to both curb the detrimental effects of overestimation bias as well as let agents discover policies that improve over those that collected the data to begin with (Levine et al., 2020). The last few years in particular have seen a wide variety of approaches proposed, making use of policy constraints (Fujimoto et al., 2019b; Zhou et al., 2021; Wu et al., 2019; Kumar et al., 2019; Kostrikov et al., 2021b; Fujimoto & Gu, 2021), conservative value estimation (Kostrikov et al., 2021a; Kumar et al., 2020), uncertainty estimation (An et al., 2021; Ghasemipour et al., 2022; Bai et al., 2022; Yang et al., 2022; Nikulin et al., 2023; Beeson & Montana, 2024) and environment modelling (Kidambi et al., 2020; Yu et al., 2021b; Argenson & Dulac-Arnold, 2020; Janner et al., 2022; Yu et al., 2020; Swazinna et al., 2021), to name just a few. Each approach comes with its own strengths and weaknesses in terms of performance, computational efficiency, ease of implementation and hyperparameter tuning.

To date, most published research in offline RL has focused on either continuous or small-scale discrete action spaces. However, many complex real-world problems can be naturally expressed in terms of *factorised action spaces*, where global actions consist of multiple distinct sub-actions, each representing a key aspect of the decision process. Examples include ride-sharing (Lin et al., 2018), recommender systems (Zhao et al., 2018), robotic assembly (Driess et al., 2020) and healthcare (Liu et al., 2020).

Formally, the factorised action space is considered the Cartesian product of a finite number of independent discrete action sets, i.e.  $\mathcal{A} = \mathcal{A}_1 \times \dots \times \mathcal{A}_N$ , where  $\mathcal{A}_i$  contains  $n_i$  (sub-)actions and  $N$  corresponds to the dimensionality of the action space. The total number of actions, often referred to as atomic actions, is thus  $\prod_{i=1}^N n_i$ , which can undergo combinatorial explosion if  $N$  and/or  $n_i$  grow large. For example, in robotics a set of actions could correspond to moving a joint up, down or not at all, i.e.  $n_i = 3$ . For machines with multiple joints the global action is the set of individual actions for each joint. The total number of possible global actions is thus  $3^N$  which grows exponentially as the number of joints increases.

In recognition of this possibility, various strategies in the online setting have been devised to preserve the effectiveness of commonly used discrete RL algorithms (Tavakoli et al., 2018; Tang & Agrawal, 2020). The concept of *value-decomposition* (Seyde et al., 2022) is particularly prominent, wherein value estimates for each sub-action space are computed independently, yet are trained to ensure that their aggregate mean converges towards a universal value. The overall effect is to reduce the total number of actions for which a value needs to be learnt from a product to a sum, making problems with factorisable actions spaces much more tractable for approaches such as Q-learning. Referring back to the robotics example, by using a decomposed approach the number of actions requiring value estimates is reduced from  $3^N$  to  $3N$ .

In this work, we undertake an initial investigation into offline RL in factorisable action spaces. Using value-decomposition, we show how a factorised approach provides several benefits over a standard atomic action representation. We conduct an extensive empirical evaluation of a number of offline approaches adapted to the factorised setting using a newly created benchmark designed to test an agent’s ability to learn complex behaviours from data of varying quality. In the spirit of advancing research in this area, we provide open access to these datasets as well as our full code base: <https://github.com/AlexBeesonWarwick/OfflineRLFactorisableActionSpaces>.

To the best of our knowledge, this investigation represents the first formative analysis of offline RL in factorisable action spaces. We believe our work helps pave the way for developments in this important domain, whilst also contributing to the growing field of offline RL more generally.

## 2 Preliminaries

### 2.1 Offline reinforcement learning

Following standard convention, we begin by defining a Markov Decision Process (MDP) with state space  $\mathcal{S}$ , action space  $\mathcal{A}$ , environment dynamics  $T(s' | s, a)$ , reward function  $R(s, a)$  and discount factor  $\gamma \in [0, 1]$  (Sutton & Barto, 2018). An agent interacts with this MDP by following a state-dependent policy  $\pi(a | s)$ ,

with the primary objective of discovering an optimal policy  $\pi^*(a | s)$  that maximises the expected discounted sum of rewards,  $\mathbb{E}_\pi [\sum_{t=0}^{\infty} \gamma^t r(s_t, a_t)]$ .

A popular approach for achieving this objective is through the use of Q-functions,  $Q^\pi(s, a)$ , which estimate the value of taking action  $a$  in state  $s$  and following policy  $\pi$  thereafter. In discrete action spaces, optimal Q-values can be obtained by repeated application of the Bellman optimality equation:

$$Q^*(s, a) = r(s, a) + \gamma \mathbb{E}_{s' \sim T} \left[ \max_{a'} Q^*(s', a') \right] .$$

These Q-values can then be used to define an implicit policy such that  $\pi(s) = \arg \max_a Q(s, a)$  i.e. the action that maximises the optimal Q-value at each state.

Given the scale and complexity of real-world tasks, Q-functions are often parameterised (Mnih et al., 2013), with learnable parameters  $\theta$  that are updated so as to minimise the following loss:

$$L(\theta) = \frac{1}{|B|} \sum_{(s, a, r, s') \sim \mathcal{B}} (Q_\theta(s, a) - y(r, s'))^2 , \quad (1)$$

where  $y(r, s') = r + \gamma \max_{a'} Q_\theta(s', a')$  is referred to as the target value, and  $B$  denotes a batch of transitions sampled uniformly at random from a replay buffer  $\mathcal{B}$ . To promote stability during training, when calculating Q-values in the target it is common to use a separate target network  $Q_{\hat{\theta}}(s', a')$  (Mnih et al., 2013; Hessel et al., 2018) with parameters  $\hat{\theta}$  updated towards  $\theta$  either via a hard reset every specified number of steps, or gradually using Polyak-averaging.

In the offline setting, an agent is prohibited from interacting with the environment and must instead learn solely from a pre-existing dataset of interactions  $\mathcal{B} = (s_b, a_b, r_b, s'_b)$ , collected from some unknown behaviour policy (or policies)  $\pi_\beta$  (Lange et al., 2012). This lack of interaction allows errors in Q-value estimates to compound and propagate during training, often resulting in a complete collapse of the learning process (Fujimoto et al., 2019b). Specifically, Q-values for out-of-distribution actions (i.e. those absent from the dataset) suffer from overestimation bias as a result of the maximisation carried out when determining target values (Thrun & Schwartz, 1993; Gaskett, 2002). The outcome is specious Q-values estimates, and policies derived from those estimates consequently being highly sub-optimal. In order to compensate for this overestimation bias, Q-values must be regularised by staying “close” to the source data (Levine et al., 2020).

## 2.2 Decoupled Q-Network

By default, standard Q-learning approaches are based on atomic representations of action spaces (Sutton & Barto, 2018). This means that, in a factorisable action space, Q-values must be determined for every possible combination of sub-actions. This potentially renders such approaches highly ineffective due to combinatorial explosion in the number of atomic actions. Recalling that the action space can be thought of as a Cartesian product, then for each  $\mathcal{A}_i$  we have that  $|\mathcal{A}_i| = n_i$ , and so the total number of atomic actions is  $\prod_{i=1}^N n_i$ . This quickly grows unwieldy as the number of sub-action spaces  $N$  and/or number of actions within each sub-action space  $n_i$  increase.

To address this issue, the Branching Dueling Q-Network (BDQ) proposed by Tavakoli et al. (2018) learns value estimates for each sub-action space independently and can be viewed as a single-agent analogue to independent Q-learning from multi-agent reinforcement learning (MARL) (Claus & Boutilier, 1998). Seyde et al. (2022) expand on this work with the introduction of the Decoupled Q-Network (DecQN), which computes value estimates in each sub-action space independently, but learns said estimates such that their mean estimates the Q-value for the combined (or global) action. Such an approach is highly reminiscent of the notion of value-decomposition used in cooperative MARL (Sunehag et al., 2017; Rashid et al., 2020b;a; Du et al., 2022), with sub-action spaces resembling individual agents.

In terms of specifics, DecQN introduces a utility function  $U_{\theta_i}^i(s, a_i)$  for each sub-action space and redefines the Q-value to be:

$$Q_\theta(s, \mathbf{a}) = \frac{1}{N} \sum_{i=1}^N U_{\theta_i}^i(s, a_i) , \quad (2)$$

where  $\mathbf{a} = (a_1, \dots, a_N)$  is the global action,  $\theta_i$  are the parameters for the  $i$ th utility function and  $\theta = \{\theta_i\}_{i=1}^N$  are the global set of parameters. The loss in Equation (1) is updated to incorporate this utility function structure:

$$L(\theta) = \frac{1}{|B|} \sum_{(s, \mathbf{a}, r, s') \sim \mathcal{B}} (Q_\theta(s, \mathbf{a}) - y(r, s'))^2, \quad (3)$$

where

$$y(r, s') = r + \frac{\gamma}{N} \sum_{i=1}^N \max_{a'_i} U_{\theta_i}^i(s', a'_i).$$

As each utility function only needs to learn about actions within its own sub-action space, this reduces the total number of actions for which a value must be learnt to  $\sum_{i=1}^N n_i$ , thus preserving the functionality of established Q-learning algorithms. Whilst there are other valid decomposition methods, in this work we focus primarily on the decomposition proposed in DecQN. In Appendix H we provide a small ablation justifying our choice.

### 3 Related Work

#### 3.1 Offline RL

Numerous approaches have been proposed to help mitigate Q-value overestimation bias in offline RL. In BCQ (Fujimoto et al., 2019b), this is achieved by cloning a behaviour policy and using generated actions to form the basis of a policy which is then optimally perturbed by a separate network. BEAR (Kumar et al., 2019), BRAC (Wu et al., 2019) and Fisher-BRC (Kostrikov et al., 2021a) also make use of cloned behaviour policies, but instead use them to minimise divergence metrics between the learned and cloned policy. One-step RL (Brandfonbrener et al., 2021; Gulcehre et al., 2020a) explores the idea of combining fitted Q-evaluation and various policy improvement methods to learn policies without having to query actions outside the data. This is expanded upon in Implicit Q-learning (IQL) (Kostrikov et al., 2021b), which substitutes fitted Q-evaluation with expectile regression. TD3-BC (Fujimoto & Gu, 2021) adapts TD3 (Fujimoto et al., 2018) to the offline setting by directly incorporating behavioural cloning into policy updates, with TD3-BC-N/SAC-BC-N (Beeson & Montana, 2024) employing ensembles of Q-functions for uncertainty estimation to alleviate issues relating to overly restrictive constraints as well as computational burden present in other ensembles based approaches such as SAC-N & EDAC (An et al., 2021), MSG (Ghasemipour et al., 2022), PBRL (Bai et al., 2022) and RORL (Yang et al., 2022).

In the majority of cases the focus is on continuous action spaces, and whilst there have been adaptations and implementations in discrete action spaces (Fujimoto et al., 2019a; Gu et al., 2022), these tend to only consider a small number of (atomic) actions. This is also reflected in benchmark datasets such as D4RL (Fu et al., 2020) and RL Unplugged (Gulcehre et al., 2020b). Our focus is on the relatively unexplored area of factorisable discrete action spaces.

#### 3.2 Action decomposition

Reinforcement learning algorithms have been extensively studied in scenarios involving large, discrete action spaces. In order to overcome the challenges inherent in such scenarios, numerous approaches have been put forward based on action sub-sampling (Van de Wiele et al., 2020; Hubert et al., 2021), action embedding (Dulac-Arnold et al., 2015; Gu et al., 2022) and curriculum learning (Farquhar et al., 2020). However, such approaches are tailored to handle action spaces comprising numerous atomic actions, and do not inherently tackle the complexities nor utilise the structure posed by factorisable actions.

For factorisable action spaces various methods have been proposed, such as learning about sub-actions independently via value-based (Sharma et al., 2017; Tavakoli et al., 2018) or policy gradient methods (Tang & Agrawal, 2020; Seyde et al., 2021). Others have also framed the problem of action selection in factorisable action spaces as a sequence prediction problem, where the sequence consists of the individual sub-actions (Metz et al., 2017; Pierrot et al., 2021; Chebotar et al., 2023).

There exists a strong connection between factorisable action spaces and MARL, where the selection of a sub-action can be thought of as an individual agent choosing its action in a multi-agent setting. Value-decomposition has been shown to be an effective approach in MARL (Sunehag et al., 2017; Rashid et al., 2020b;a; Du et al., 2022), utilising the *centralised training with decentralised execution* paradigm (Kraemer & Banerjee, 2016), which allows agents to act independently but learn collectively. DecQN (Seyde et al., 2022) and REValueD (Ireland & Montana, 2023) have subsequently shown such ideas can be used with factorised action spaces in single-agent reinforcement learning, demonstrating strong performance on a range of tasks that vary in complexity. Theoretical analysis of errors in Q-value estimates has also been conducted in efforts to stabilise training (Ireland & Montana, 2023; Thrun & Schwartz, 1993).

In this work, we focus on adapting DecQN to the offline setting by incorporating existing offline techniques. Whilst prior work has explored offline RL with value decomposition (Tang et al., 2022), this was limited to specific low-dimensional healthcare applications using only BCQ. Furthermore, accurately evaluating performance in such domains is notorious challenging (Gottesman et al., 2018). In contrast, we systematically study multiple offline methods using low and high-dimensional factorised action spaces across a suite of benchmark tasks.

## 4 A case for factorisation and decomposition in offline RL

As mentioned in Section 2.2, value-decomposition provides a mechanism for overcoming challenges in standard Q-learning arising from exponential growth in the number of atomic actions. This can be seen most clearly from a practical standpoint, in which the number of actions that require value estimation is significantly reduced when moving from atomic to factorised action representation (see Table 4 for example). However, even in cases where atomic representation remains feasible, there are still several benefits to a factorised and decomposed approach which are particularly salient in the offline case.

To see this, we begin by making the assumption that Q-value estimates from function approximation  $Q_\theta(s, \mathbf{a})$  carry some noise  $\epsilon(s, \mathbf{a})$  (Thrun & Schwartz, 1993). This noise can be modelled as a random variable, such that  $Q_\theta(s, \mathbf{a}) = Q^\pi(s, \mathbf{a}) + \epsilon(s, \mathbf{a})$ , where  $Q^\pi(s, \mathbf{a})$  is the true but unknown Q-value for policy  $\pi$ . For clarity, we emphasise this noise stems from the fact the function approximator cannot represent the true Q-function precisely, as opposed to other factors such as sampling variation. As noted by Thrun & Schwartz (1993), the expected overestimation will be maximal if all actions for a particular state share the same target Q-value. Incorporating all these ideas into the Q-learning framework we can define the target difference as

$$\begin{aligned} Z_{s'}^{dqN} &= \gamma \left( \max_{\mathbf{a}'} Q_\theta(s', \mathbf{a}') - \max_{\mathbf{a}'} Q^\pi(s', \mathbf{a}') \right); \\ &= \gamma \left( \max_{\mathbf{a}'} \epsilon(s', \mathbf{a}') \right). \end{aligned} \tag{4}$$

To incorporate decomposition as defined in Equation 2, the above can be extended to utility value estimates (Ireland & Montana, 2023). Now we assume that utility value estimates from function approximation  $U_{\theta_i}^i(s, a_i)$  carry some noise  $\epsilon^i(s, a_i)$ , such that  $U_{\theta_i}^i(s, a_i) = U_i^{\pi_i}(s, a_i) + \epsilon^i(s, a_i)$ , where  $U_i^{\pi_i}(s, a_i)$  is the true but unknown utility value for policy  $\pi_i$ . Under the additional assumption that true Q-values decompose in the same way as Equation 2 (Ireland & Montana, 2023), we can define the target difference under DecQN as

$$\begin{aligned} Z_{s'}^{dec} &= \gamma \left( \frac{1}{N} \sum_{i=1}^N \max_{a'_i} U_{\theta_i}^i(s', a'_i) - \frac{1}{N} \sum_{i=1}^N \max_{a'_i} U_i^{\pi_i}(s', a'_i) \right); \\ &= \gamma \left( \frac{1}{N} \sum_{i=1}^N \max_{a'_i} \epsilon^i(s', a'_i) \right). \end{aligned} \tag{5}$$

We can expand on each of these by considering the composition of global actions in terms of in-distribution and out-of-distribution sub-actions. When actions are represented atomically, a particular global action  $\mathbf{a}$  is either in-distribution or out-of-distribution, depending on its presence (or absence) in the dataset. However, under a factorised representation a global action  $\mathbf{a}$  is composed of sub-actions  $(a_1, a_2, \dots, a_N)$ , and

hence it is the individual sub-actions that are either in-distribution or out-of-distribution. This means that global actions which are out-of-distribution can contain individual sub-actions that are in-distribution when factorised, as illustrated in Figure 1.

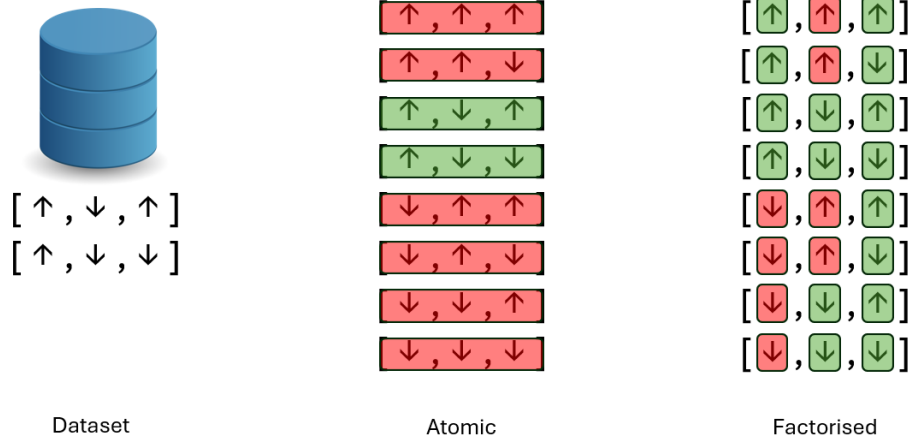


Figure 1: In this simple example there are  $N = 3$  sub-action dimensions, each with two actions  $\{\uparrow, \downarrow\}$ . In-distribution and out-of-distribution actions/sub-actions are highlighted in green and red, respectively. For a particular state, the dataset contains two global actions. Under atomic representation only actions which match those in the dataset are in-distribution. Under factorised representation, individual sub-actions which match those in the dataset are in-distribution. Atomic actions that are out-of-distribution can contain sub-actions that are in-distribution when factorised.

In recognition of this property, we now differentiate between errors in value estimates for in-distribution and out-of-distribution actions. Using the precept that out-of-distribution errors are inherently larger than in-distribution errors, let  $\epsilon^{in}(s, \mathbf{a}^{in})$ ,  $\epsilon^{in}(s, a_i^{in})$  and  $\epsilon^{out}(s, \mathbf{a}^{out})$ ,  $\epsilon^{out}(s, a_i^{out})$  denote the errors in Q-value/utility value estimates for in-distribution  $\mathbf{a}^{in}, a_i^{in}$  and out-of-distribution  $\mathbf{a}^{out}, a_i^{out}$  global/sub-actions, respectively. The target difference for DQN and DecQN under this framework become

$$Z_s^{dqn} = \gamma \left( \max_{\mathbf{a}^{in}} \{ \max_{a_i^{in}} \epsilon^{in}(s, a_i^{in}), \max_{\mathbf{a}^{out}} \epsilon^{out}(s, \mathbf{a}^{out}) \} \right); \quad (6)$$

$$Z_s^{dec} = \gamma \left( \frac{1}{N} \sum_{i=1}^N \max_{a_i^{in}} \{ \max_{a_i^{in}} \epsilon^{in}(s, a_i^{in}), \max_{a_i^{out}} \epsilon^{in}(s, a_i^{out}) \} \right); \quad (7)$$

where we have dropped the prime symbol to avoid overloading notation. Now the maximum is taken over a pooled sample of errors from two different distributions as opposed to a single distribution in Equations 4 and 5.

The important point to note here is that in both cases the target difference is now dependent on the relative number of in-distribution and out-of-distribution actions/sub-actions. Since errors from in-distribution actions/sub-actions are smaller than for out-of-distribution, the larger the number of in-distribution actions/sub-actions the smaller the expected overestimation. In the offline setting, we cannot change the status of an action/sub-action from out-of-distribution to in-distribution through environment interaction. However, as noted above, we can improve the coverage of sub-actions through factorisation. Hence, for the same dataset we can potentially reduce overestimation bias moving from an atomic action representation to a factorised one. This is particularly beneficial in the offline setting where issues relating to overestimation bias primarily stem from Q-value estimates for out-of-distribution actions.

In Appendix A we provide support for the ideas presented in this Section by empirically assessing the properties of target differences under DQN and DecQN for uniformly distributed noise.

#### 4.1 Appropriateness of decomposition in the offline setting

The previous analysis assumes both the approximate and true Q-function decompose as per Equation 2. In general, such a decomposition is unlikely to perfectly capture intricate aspects of an environment. Nonetheless, such a decomposition has been shown to be effective in the online setting (Ireland & Montana, 2023), and we note that under the following conditions it can act as reasonable approximation, allowing us to retain the aforementioned benefits of improved action/sub-action coverage.

##### Weak inter-action dependence

If the effect of each sub-action on the reward or the next state is relatively independent of other sub-actions, the decomposition is more likely to hold. This condition aligns with the notion of action independence in factored MDPs, where the transition dynamics can be decomposed across sub-actions with minimal cross-interaction. In such cases, the global Q-function can be closely approximated by a sum of local utility functions. This can also be tied to work on factored MDPs (Kearns & Koller, 1999), where state and action space factorisations allow efficient policy computation under certain independence assumptions. Of course, true independence of sub-actions in an MDP might be rare in practice, especially when complex dependencies exist between sub-actions (e.g., coordination between multiple control signals). Still, if dependencies are weak, factorisation can be a good approximation. It’s also important to emphasise this only holds under certain independence assumptions that often arise in factored MDPs but may not be present in more general MDPs.

##### Approximately factorisable reward structures

The decomposition is also reasonable when the reward function can be roughly expressed as a sum over individual sub-actions, which is again related to the literature on reward decomposition in factored MDPs (Guestrin et al., 2003). In scenarios where the global reward is a combination of local rewards for each sub-action, factorisation provides a natural approximation for value function decomposition. The caveat is that, in many real-world scenarios, rewards are not perfectly factorisable, and there could be interactions between sub-actions that affect the total reward in a non-linear way. Thus, such an approximation may break down in complex environments with highly entangled reward structures. However for certain structured problems, especially those with weak interactions, this approximation can still yield practical benefits.

#### 4.2 Limitations of decomposition in the offline setting

It is important to highlight that while the benefits of a decomposed approach in the offline setting can offset shortcomings in modelling inter-action dependence, there are limitations. For settings where sub-actions are strongly dependent, the value of a global action may differ significantly to the average value of its constituent utilities, introducing errors that are no longer outweighed by the reduction in overestimation bias from factorisation/decomposition.

As an example, consider the phenomena of pharmacodynamic and pharmacokinetic drug interactions. In a healthcare setting, it is common for patients to receive multiple drugs as part of a treatment regime, with dosages and routes adjusted based on how the patient responds. If our goal is to train an offline RL to optimise treatment regimes, the level of interaction (positive or negative) between drugs/dosages/route would dictate whether a factorised/decomposed approach would be appropriate. If there is little to no interaction then a factorised/decomposed approach could be effective. If there are strong interactions a factorised/decomposed approach would likely fail to capture these important characteristics and hence be much less effective.

## 5 Algorithms

In this Section we introduce several approaches for adapting DecQN to the offline setting based on existing techniques. We focus on methods that offer distinct takes on combatting overestimation bias, namely, policy constraints, conservative value estimation, implicit Q-learning and one-step RL. In each case, attention shifts

from Q-values to utility values, with regularisation performed at the sub-action level. The full procedure for each algorithm can be found in Appendix B

### 5.1 DecQN-BCQ

Batch Constrained Q-learning (BCQ) (Fujimoto et al., 2019b;a) is a policy constraint approach to offline RL. To compensate for overestimation bias for out-of-distribution actions, a cloned behaviour policy  $\pi_\phi$  is used to restrict the actions available for target Q-values estimates, such that their probability under the behaviour policy meets a relative threshold  $\tau$ . This can be adapted and integrated into DecQN by cloning separate behaviour policies  $\pi_{\phi_i}^i$  for each sub-action dimension and restricting respective sub-actions available for corresponding target utility value-estimates. The target value from Equation (3) becomes:

$$y(r, s') = r + \frac{\gamma}{N} \sum_{i=1}^N \max_{a'_i : \rho^i(a'_i) \geq \tau} U_{\theta_i}^i(s', a'_i),$$

where  $\rho^i(a'_i) = \pi_{\phi_i}^i(a'_i | s') / \max_{\hat{a}'_i} \pi_{\phi_i}^i(\hat{a}'_i | s')$  is the relative probability of sub-action  $a'_i$  under policy  $\pi_{\phi_i}^i$ . Each cloned behaviour policy is trained via supervised learning with  $\phi = \{\phi\}_{i=1}^N$ . The full procedure can be found in Algorithm 1.

### 5.2 DecQN-CQL

Conservative Q-learning (CQL) (Kumar et al., 2020) attempts to combat overestimation bias by targeting Q-values directly. The loss in Equation (1) is augmented with a term that “pushes-up” on Q-value estimates for actions present in the dataset and “pushes-down” for all others. Under one particular variant this additional loss takes the form

$$L_{CQL}(\theta) = \frac{\alpha}{|B|} \sum_{s, a \sim \mathcal{B}} \left[ \log \sum_{a_i \in A} \exp(Q_\theta(s, a_i)) - Q_\theta(s, a) \right];$$

which equates to performing behavioural cloning using log-likelihood when the policy is a softmax over Q-values (Luo et al., 2023) since

$$\begin{aligned} \frac{\alpha}{|B|} \sum_{s, a \sim \mathcal{B}} \left[ \log \sum_{a_i \in A} \exp(Q_\theta(s, a_i)) - Q_\theta(s, a) \right] &= -\frac{\alpha}{|B|} \sum_{s, a \sim \mathcal{B}} \left[ \log \frac{\exp(Q_\theta(s, a))}{\sum_{a_i \in A} \exp(Q_\theta(s, a_i))} \right] \\ &= -\frac{\alpha}{|B|} \sum_{s, a \sim \mathcal{B}} \log \pi_\theta(a | s). \end{aligned}$$

This variant can be adapted and incorporated into DecQN by “pushing-up” on utility value estimates for sub-actions present in data and “pushing-down” for all others. The additional loss under DecQN becomes:

$$L_{CQL}(\theta) = \frac{\alpha}{|B|} \sum_{s, \mathbf{a} \sim \mathcal{B}} \frac{1}{N} \sum_{i=1}^N \left[ \log \sum_{a_j \in A_i} \exp(U_{\theta_i}^i(s, a_j)) - U_{\theta_i}^i(s, a_i) \right]; \quad (8)$$

where  $a_j$  denotes the  $j$ th sub-action within the  $i$ th sub-action space, and  $\alpha$  is a hyperparameter that controls the overall level of conservatism. The full procedure can be found in Algorithm 2.

Note this particular implementation does not directly substitute the decomposition as per Equation 2 into the original CQL loss. Instead, we apply CQL directly to the utility values and then take the mean across sub-action dimensions. This way we avoid having to estimate Q-values for all atomic actions (which a direct substitution would require) and we retain the equivalence to behavioural cloning at the sub-action level as

$$\begin{aligned} \frac{\alpha}{|B|} \sum_{s, \mathbf{a} \sim \mathcal{B}} \frac{1}{N} \sum_{i=1}^N \log \sum_{a_j \in A_i} \exp(U_{\theta_i}^i(s, a_j)) - U_{\theta_i}^i(s, a_i) &= -\frac{\alpha}{|B|} \sum_{s, \mathbf{a} \sim \mathcal{B}} \frac{1}{N} \sum_{i=1}^N \left[ \log \frac{\exp(U_{\theta_i}^i(s, a_i))}{\sum_{a_j \in A_i} \exp(U_{\theta_i}^i(s, a_j))} \right] \\ &= -\frac{\alpha}{|B|} \sum_{s, \mathbf{a} \sim \mathcal{B}} \frac{1}{N} \sum_{i=1}^N \log \pi_{\theta_i}(a_i | s). \end{aligned}$$



### 5.3 DecQN-IQL

Implicit Q-learning (IQL) (Kostrikov et al., 2021b) addresses the challenge of overestimation bias by learning a policy without having to query actions absent from data. A state and state-action value function are trained on the data and then used to extract a policy via advantage-weighted-behavioural-cloning.

The state value function  $V_\psi(s)$  is trained via expectile regression, minimising the following loss:

$$L(\psi) = \frac{1}{|B|} \sum_{s,a \sim \mathcal{B}} [L_2^\tau(Q_\theta(s,a) - V_\psi(s))];$$

where if we denote  $u = Q_\theta(s,a) - V_\psi(s)$  then  $L_2^\tau(u) = |\tau - 1(u < 0)|u^2$  is the asymmetric least squares for the  $\tau \in (0, 1)$  expectile. The state-action value function  $Q_\theta(s,a)$  is trained using the same loss as Equation (1), with the target value now  $y(r,s') = r + \gamma V_\psi(s')$ .

The policy follows that of discrete-action advantage-weighted-behavioural-cloning (Luo et al., 2023), which adapted to the factorised setting and integrated with DecQN leads to the following:

$$\pi_i = \arg \max_{a_i} \left[ \frac{1}{\lambda} A(s, a_i) + \log \pi_{\phi_i}^i(a_i | s) \right].$$

where  $A(s, a_i) = U(s, a_i) - V(s)$  is the (factorised) advantage function,  $\pi_{\phi_i}^i(a_i | s)$  is the cloned behaviour policy for sub-action space  $i$  trained via supervised learning and  $\lambda$  is a hyperparameter controlling the balance between reinforcement learning and behavioural cloning. The full procedure can be found in Algorithm 3.

### 5.4 DecQN-OneStep

We can derive an alternative approach to IQL which removes the requirement for a separate state value function altogether. Noting that  $V(s) = \sum_a \pi(a | s) Q(s, a)$ , we can instead use the cloned behaviour policy  $\pi_\phi(a' | s')$  and state-action value function  $Q_\theta(s', a')$  to calculate the state value function  $V(s')$  instead. This can be adapted and incorporated into DecQN by replacing  $Q(s, a)$  with its decomposed form as per Equation (2) and adjusting the policy to reflect a sub-action structure. We denote this approach DecQN-OneStep as it mirrors one-step RL approaches that train state value functions using fitted Q-evaluation (Brandfonbrener et al., 2021). The full procedure can be found in Algorithm 4.

## 6 Environments and datasets

At present, there are relatively few established environments/tasks specifically designed for factorised action spaces. As such, there is an absence of benchmark datasets akin to those available for continuous or small-scale discrete action spaces such as D4RL (Fu et al., 2020) and RL Unplugged (Gulcehre et al., 2020b). In light of this, we introduce our own benchmarking suite constituting the following environments and tasks.

**Maze:** A simple maze-based environment first introduced by Chandak et al. (2019) in which an agent is tasked with reaching a target goal location. The state space is continuous and comprises the agent’s current location in the maze. The agent’s movement is controlled by a series of  $N$  actuators which can be turned on or off, with each actuator corresponding to a unit of movement in a single direction, with actuator  $i$  applying force at an angle of  $2\pi i/N$  when activated, as illustrated in Figure 2. The net outcome of actuator selection is the vectoral summation of the movements associated with each actuator. The agent’s action space comprises the unique combination of the set of binary decision with respect to each actuator, leading to to an action set that is exponential in the number of actuators such that  $|A| = 2^N$ .

**DeepMind control suite:** A discretised variant of the DeepMind control suite (Tunyasuvunakool et al., 2020), as previously adopted by Seyde et al. (2022); Ireland & Montana (2023). This suite contains a variety of environments and tasks, which although originally designed for continuous control can easily be repurposed for a discrete factorisable setting by discretising individual sub-action spaces, thus creating scenarios that can vary in size and complexity. This discretisation process involves selecting a subset of actions from the original continuous space which then become the discrete actions available to the agent. For example, a

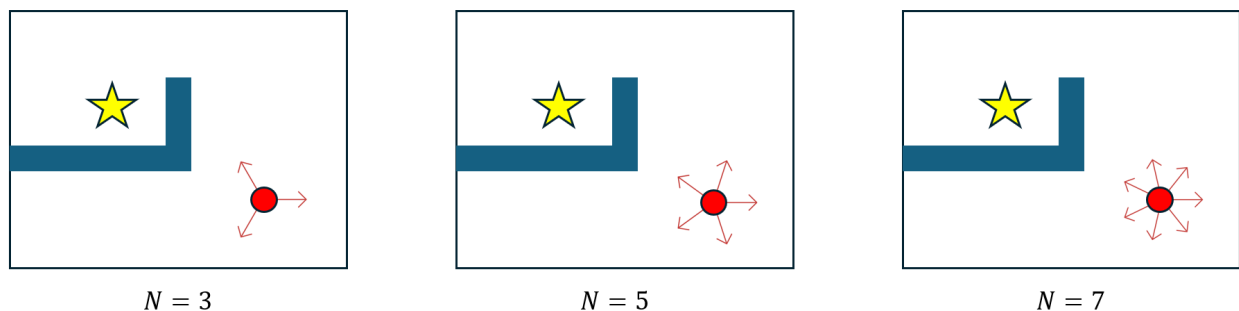


Figure 2: Examples of maze environment with different numbers of actuators. The star represents the target goal location, the red dot the agent and the arrows the actuators. Adapted from original Figure in Chandak et al. (2019).

continuous action in the range  $[-1, 1]$  can be discretised into three discrete actions corresponding to the subset  $\{-1, 0, 1\}$ . We emphasise this discretisation procedure happens prior to data collection, i.e. we are not discretising an existing continuous action dataset.

The choice of environments and tasks reflects our desire to provide a benchmark that is amenable to factorisation and decomposition to varying degrees, allowing us to investigate the limits of the assumptions and conditions outlined in Section 4. For the maze task, the action set is inherently factorisable for each actuator, yet reaching the target goal requires at least some coordination across actuators. For the DeepMind control suite, each sub-action corresponds to the adjustment of a specific robotic component (e.g. a joint), but to successfully complete the task may require alignment across some or all components. We emphasise that our benchmark is not intended to cover large scale discrete action settings for which factorisation/decomposition is not feasible.

For the datasets themselves, we follow a similar procedure to D4RL. Using DecQN/REValueD, we train agents to “expert” and “medium” levels of performance and then collect transitions from the resulting policies. Here, we define “expert” to be the peak performance achieved by DecQN/REValueD and “medium” to be approximately 1/3rd the performance of the “expert”. We create a third dataset “medium-expert” by combining transitions from these two sources and a fourth “random-medium-expert” containing transitions constituting 45% random and medium transitions and 10% expert. Each of these datasets presents a specific challenge to agents, namely the ability to learn from optimal or sub-optimal data (“expert” and “medium”, respectively) as well as data that contains a mixture (“medium-expert” and “random-medium-expert”). More details on this training and data collection procedure are provided in Appendix C.

## 7 Experimental evaluation

We train agents using DecQN, DecQN-BCQ, DecQN-CQL, DecQN-IQL and DecQN-OneStep on our benchmark datasets and evaluate their performance in the simulated environment. We also train and evaluate agents using a factorised equivalent of behavioural cloning to provide a supervised learning baseline. Performance is measured in terms of normalised score, where  $score_{norm} = 100 \times \frac{score - score_{random}}{score_{expert} - score_{random}}$  with 0 representing a random policy and 100 the “expert” policy from the fully trained agent. We repeat experiments across five random seeds, reporting results as mean normalised scores  $\pm$  one standard error across seeds. For each set of experiments we provide visual summaries with tabulated results available in Appendix G for completeness. Full implementation details, including network architectures, hyperparameters and training procedures are provided in Appendix D

### 7.1 Case study: DQN-CQL vs DecQN-CQL

Before considering the full benchmark, we conduct a case study directly comparing DQN-CQL (i.e. CQL under atomic action representation) and DecQN-CQL. Following the procedure outlined in Section 6 we construct a “medium-expert” dataset for the “cheetah-run” task for number of bins  $n_i \in \{3, 4, 5, 6\}$  and compare the performance of resulting policies and overall computation time. In addition we construct a “random-medium-expert” dataset for the Maze task with  $N = 15$  actuators for varying dataset sizes and compare performance and computational resources.

In Figure 3 we see that as the number of sub-actions  $n_i$  increases for the “cheetah-run” task, DQN-CQL exhibits a notable decline in performance, whereas DecQN-CQL performance declines only marginally. We also see a dramatic increase in computation time for DQN-CQL whereas DecQN-CQL remains roughly constant. For DQN-CQL, these outcomes are symptomatic of the combinatorial explosion in the number of actions requiring value-estimation and the associated number of out-of-distribution global actions. These issues are less prevalent in DecQN-CQL due to its factorised and decomposed formulation. To provide further insights, in Appendix F we also examine the evolution of Q-values errors during training, finding these errors are consistently lower for DecQN-CQL than DQN-CQL, aligning with each algorithm’s respective performance.

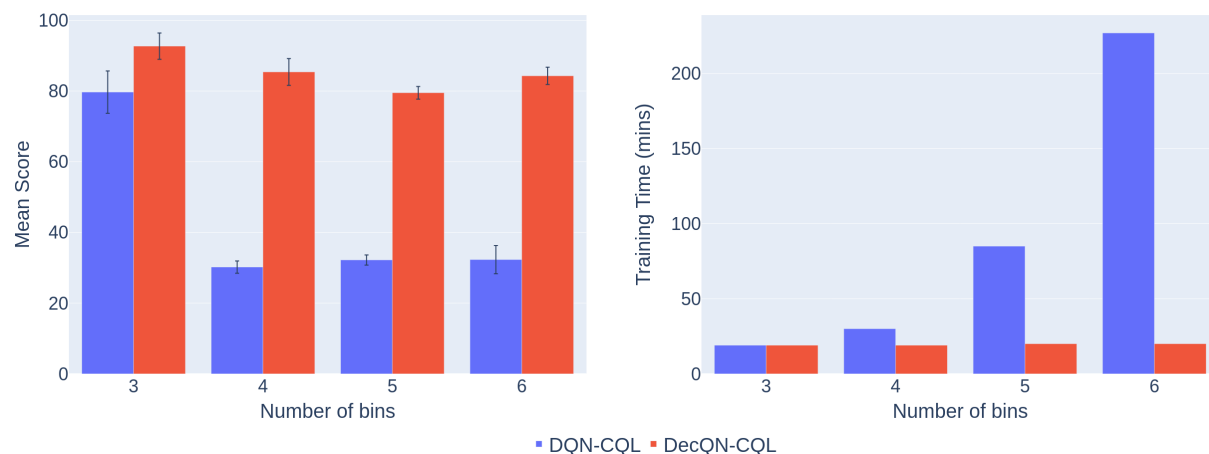


Figure 3: Comparisons of performance (left) and computation time (right) for DQN-CQL and DecQN-CQL on the “cheetah-run-medium-expert” dataset for varying numbers of bins. As the number of bins increases, DQN-CQL suffers notable drops in performance and increases in computation time, whereas DecQN-CQL is relatively resilient in both areas.

In Figure 4 we see that as the size of the dataset for the Maze task decreases, DQN-CQL exhibits a more notable decline in performance than DecQN-CQL. This is particularly the case when the size of the dataset is very small, i.e.  $\leq 250$  transitions. These outcomes support our notion of better action coverage under a factorised representation than atomic, where value estimates for observed sub-actions can be combined to provide better estimates of global actions that haven’t been observed. Furthermore, this is achieved through much more computationally efficient means, with training time for DecQN-CQL being over 8 times faster than DQN-CQL (4mins vs 34mins, respectively) and GPU usage 7 times less (246MBs vs 1728MBs, respectively).

### 7.2 Maze

In Figure 5 we summarise each algorithm’s performance for varying numbers of actuators and dataset composition. Across the board, we see that DecQN sans offline modification results in policies that perform no

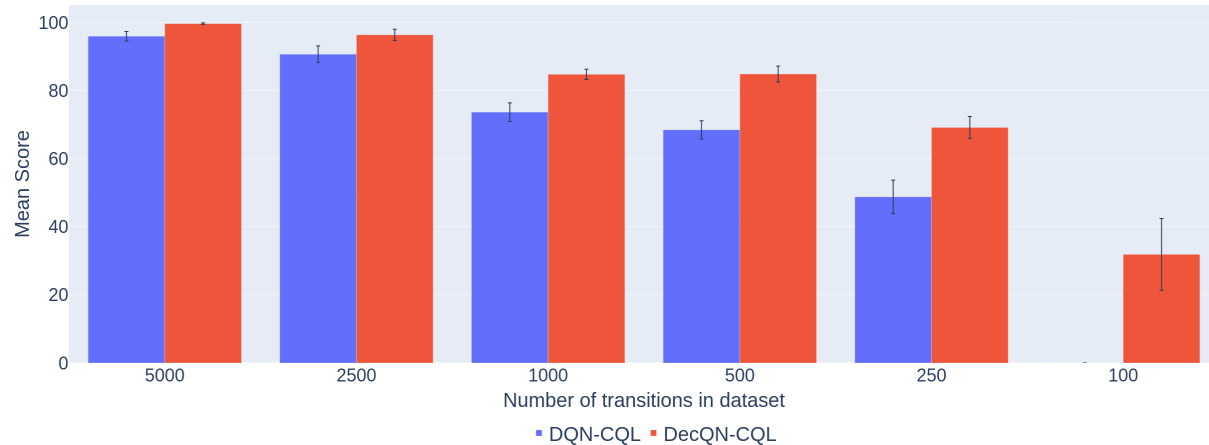


Figure 4: Comparison of performance for DQN-CQL and DecQN-CQL on the Maze task with  $N = 15$  actuators with “random-medium-expert” datasets of varying size. As the number of transitions in the dataset decreases, DQN-CQL suffers more notable drops in performance in comparison to DecQN-CQL.

better than random, a direct consequence of aforementioned issues relating to overestimation bias. With the exception of DecQN-BCQ for “random-medium-expert” datasets, all offline methods match or outperform behavioural cloning, particularly for datasets constituting high levels of sub-optimal trajectories (“medium” and “random-medium-expert”). For datasets containing a proportion of expert trajectories (“medium-expert” and “random-medium-expert”) all offline methods (again with the exception of DecQN-BCQ) are able to extract expert or near-expert level policies.<sup>1</sup> Comparing just the offline RL methods, we see there is little to separate them for “expert” and “medium-expert” regardless of the number of actuators. For “medium” and “random-medium-expert”, DecQN-CQL exhibits an edge when the number of actuators is between 3 and 12, but for “medium” this is lost to DecQN-IQL/OneStep as the number of actuators increases to 15.

### 7.3 DeepMind Control Suite

In Figure 6 we summarise each algorithm’s performance across the full range of tasks, setting the number of sub-actions  $n_i = 3$ . This necessitates the use of value-decomposition for all but the most simple tasks, as highlighted in Table 4. In general we see that all offline RL methods outperform behavioural cloning across all environments/tasks and datasets, with the exception of DecQN-BCQ for “random-medium-expert” datasets which performs quite poorly. In terms of offline methods specifically, in general DecQN-CQL has a slight edge over others for lower dimensional tasks such as “finger-spin”, “fish-swim” and “cheetah-run”, whilst DecQN-IQL/OneStep have the edge for higher-dimensional tasks such as “humanoid-stand” and “dog-trot”. For “medium-expert” datasets we see in most cases all offline methods are able to learn expert or near-expert level policies. Extracting optimal behaviour from “random-medium-expert” datasets proves significantly more challenging than in the Maze environment, likely a result of these environments/tasks being much more complex coupled with datasets being both highly variable and constituting relatively few expert trajectories in relation to trajectory length.

#### 7.3.1 Increasing the number of sub-actions

To help provide insights into the ability of our chosen offline methods to scale to larger and larger action spaces, we increase the number of sub-actions within each sub-action space and repeat our experiments. We

<sup>1</sup>For ‘medium-expert’ datasets we note that BC is able to achieve expert-level performance on par with our offline methods. We attribute this to the relative simplicity of the environment coupled with a reasonable level of expert-level trajectories within the datasets. Such outcomes are not replicated in DMC tasks which are more complex and have much longer trajectories

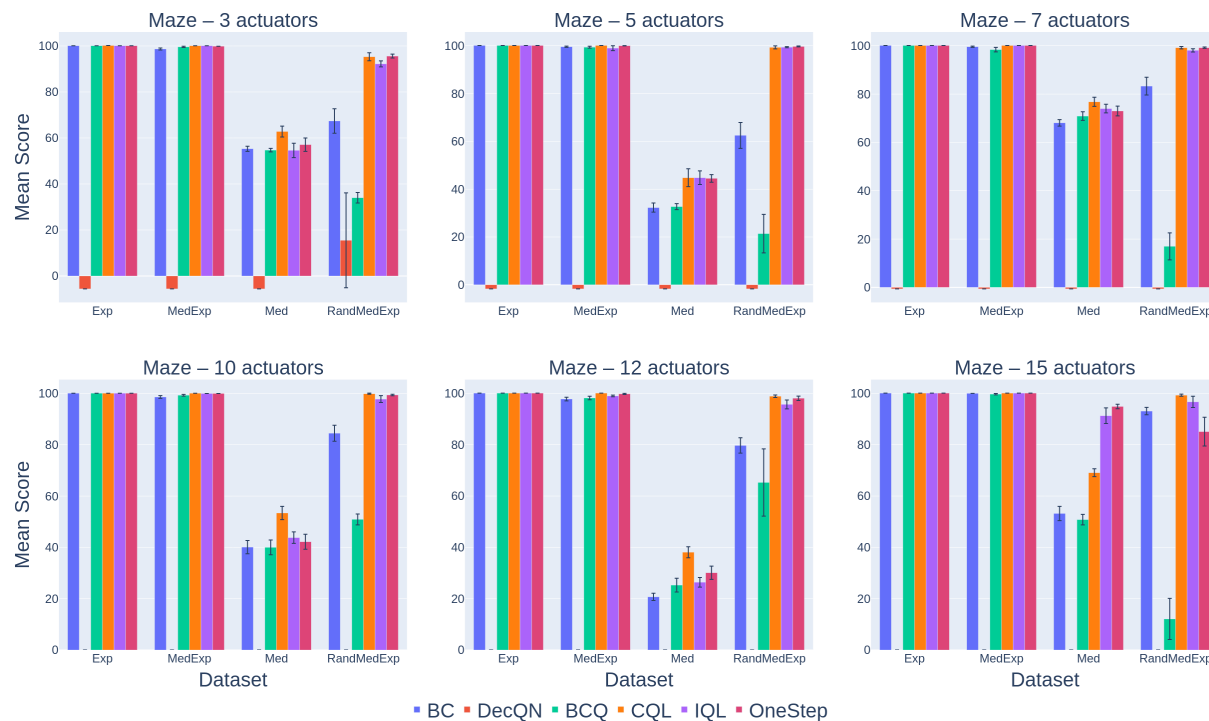


Figure 5: Performance comparison on maze task for varying numbers of actuators. For presentation purposes the prefix “DecQN-” has been omitted for each offline method. In general, all offline RL methods improve over behavioural cloning, with the exception of DecQN-BCQ for “random-medium-expert” datasets. DecQN without any offline modification performs poorly across the board.

focus in particularly on the dog-trot environment since this is by far the largest in terms of actions. We collect datasets following the same procedure outlined in Section 6 for number of bins  $n_i \in \{10, 30, 50, 75, 100\}$ .

We summarise results in Figure 7. In general, we see that our chosen offline methods are robust to increases in the number of bins, continuing to outperform behavioural cloning (with the same exception for DecQN-BCQ on “random-medium-expert”) and extract near-expert policies from “medium-expert” datasets, with DecQN-IQL/-OneStep maintaining their edge over DecQN-CQL. For “random-medium-expert” datasets we start to notice a decline in performance as we approach the upper end of our number of bins, most noticeably when  $n = 100$ . This is likely a consequence of more and more actions exacerbating the difficulties in obtaining good policies from highly variable and largely sub-optimal data.

## 8 Discussion and conclusion

In this work, we have conducted a formative investigation of offline reinforcement in factorisable action spaces. Through empirical evaluation, we have shown how a factorised and decomposed approach offers numerous benefits over standard/atomic approaches. Using a bespoke benchmark we have undertaken an extensive empirical evaluation of several offline RL approaches adapted to this factorised and decomposed framework, providing insights into each approach’s ability to learn tasks of varying complexity from datasets of differing size and quality.

In general, our empirical evaluation demonstrates our chosen offline methods adapt well to the factorised setting when combined with value-decomposition in the form of DecQN. With one exception, all approaches are consistently able to outperform behavioural cloning regardless of data quality, and where datasets contain

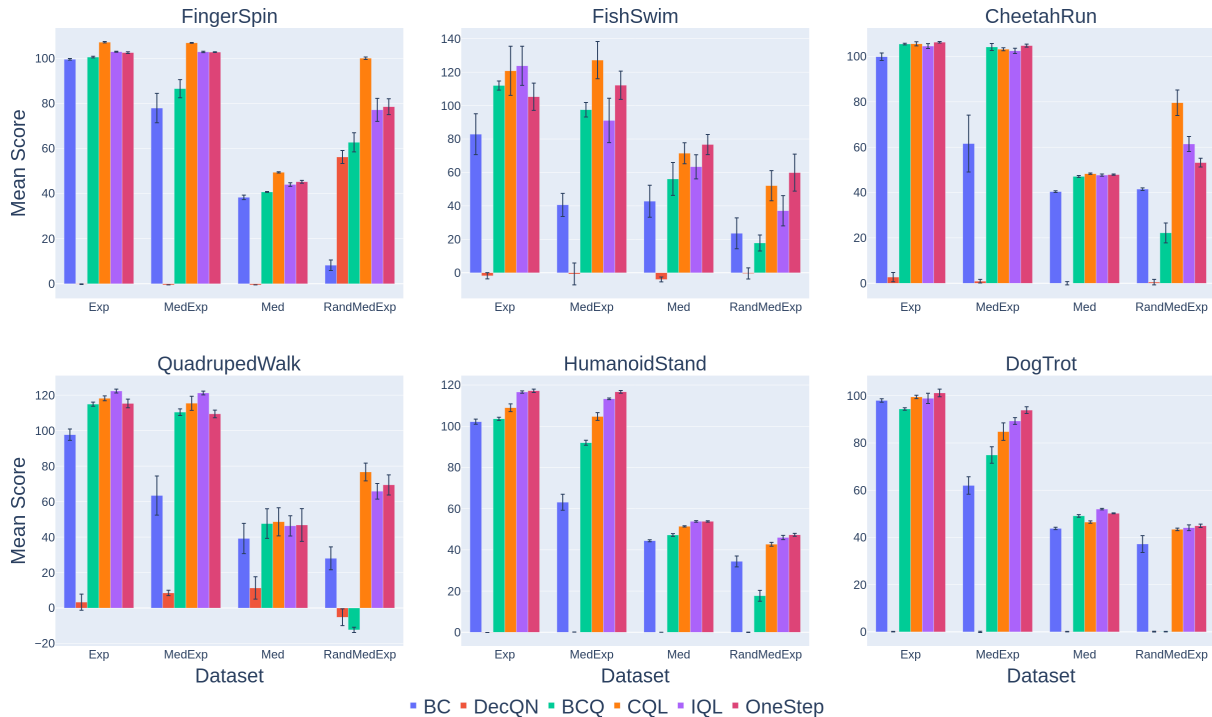


Figure 6: Performance comparison across benchmark for  $n_i = 3$ . In general, all offline RL methods improve over behavioural cloning, with the exception of DecQN-BCQ for “random-medium-expert” datasets. DecQN without any offline modification performs poorly across all environments/tasks.

sufficient levels of high-quality trajectories (i.e. “expert” and “medium-expert”), obtain expert/near-expert policies. There is however notable room for improvement for datasets with a scarcity of high-quality trajectories for more complex tasks (i.e. “medium” and “random-medium-expert” for DMC environments/tasks).

Our initial investigation opens up numerous other possibilities for future research. One of these is the development of techniques for automatically tuning hyperparameters during training, which at present are not environment/task agnostic (only data quality). In addition, as with continuous counterparts, performance can be enhanced by allowing hyperparameters to vary for each dataset (see Appendix D.1 for examples). Off-policy evaluation could also prove beneficial here (Rebello et al., 2023), providing assurances on the quality of a policy prior to deployment.

For DecQN-BCQ/-IQL/-OneStep, alternative approaches to modelling the behaviour policy  $\pi_\phi$  may help improve performance for more challenging datasets. Incorporating other methods outlined in Section 3 may also prove beneficial. For example, the use of ensembles for capturing uncertainty in value estimates has been shown to perform well in combination with behavioural cloning in continuous action spaces (Beeson & Montana, 2024), and is a relatively straightforward extension to the approaches we consider here.

Whilst DecQN offers a simple, effective and computationally efficient foundation for offline RL in factorisable action spaces, we note there are some inherent assumptions and limitations to the value decomposition approach that warrant further investigation, as noted in Sections 4.1 - 4.2. In particular, the efficacy of DecQN relies on individual sub-action optimisation leading to globally optimal joint policies. However, for tasks with sparse rewards or complex sub-action dependencies, individually learned sub-policies may fail to properly compose into a coherent overall policy. For example, in assembly tasks, separately learned pick, place, and connect skills could lead to conflicting behaviors when combined. Additional research into modeling sub-action interactions during decomposition could help overcome this limitation.

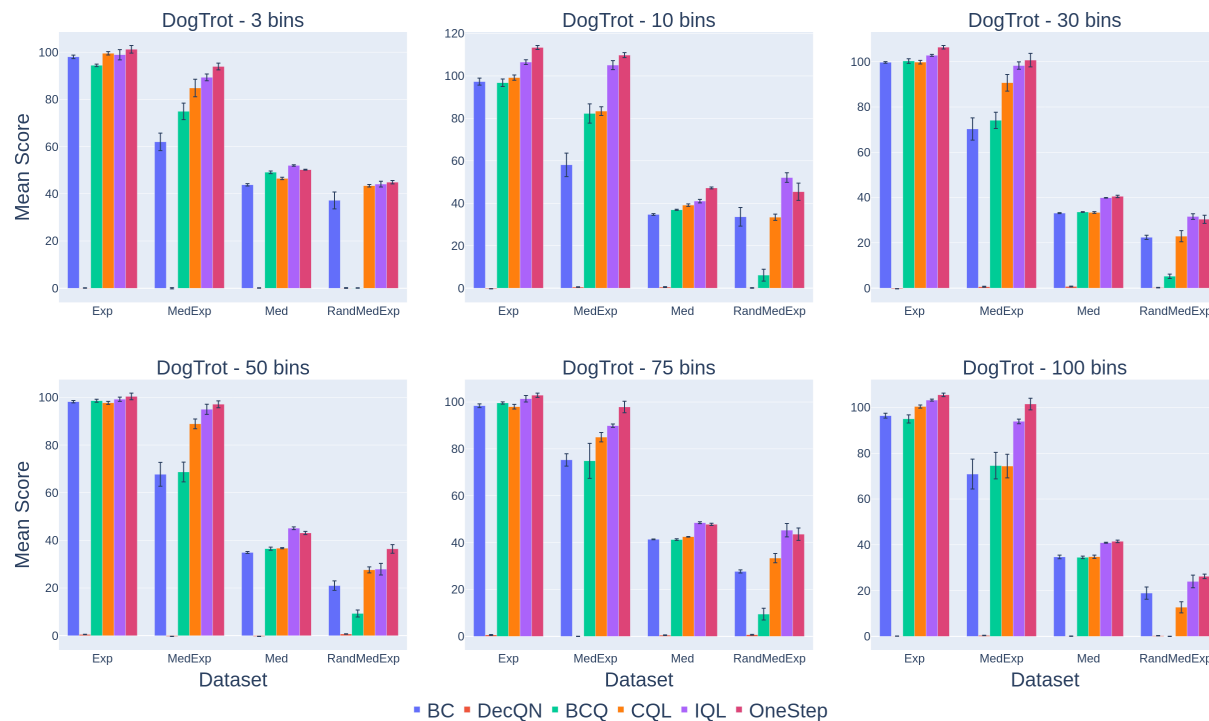


Figure 7: Performance comparison for dog-trot for  $n_i \in \{3, 10, 30, 50, 75, 100\}$ . Each approach is reasonably resilient to increases in the number of bins, although for “random-medium-expert” datasets extracting a good policy appears to become more challenging as  $n$  gets very large.

Our technical analysis can also potentially be expanded upon by considering the degree to which a global action is in-distribution/out-of-distribution. Implicitly we have assumed errors from actions/sub-actions are the same regardless of how many times a global action appears in the dataset. Relaxing this assumption could help refine the analysis, particularly for stochastic environments where multiple interactions can provide additional information. Further refinement may also be possible by addressing assumptions on Q-function decomposition (as highlighted in Sections 4.1 - 4.2), particularly the true Q-function which may not necessarily decompose in the same manner as DecQN.

Beyond algorithmic considerations, there are also several avenues for exploration in relation to the environments and tasks themselves. While our benchmark provides a broad range of diverse tasks, any particular one is deficient in either complexity (i.e. maze) or being naturally factorisable (i.e. DMC). Thus, the creation of bespoke environments that help bridge this gap would be a valuable contribution, providing more realistic scenarios to evaluate against. In addition, future work could investigate more nuanced aspects relating to action space discretisation, such as variable numbers of bins, non-even spacing between actions, clustering of actions and masked actions.

We hope our work underscores the unique setting and challenges of conducting offline RL in factorisable action spaces and paves the way for future research by providing an accessible and solid foundation from which to build upon.

**Acknowledgments** AB acknowledges support from University of Warwick and University of Birmingham NHS Foundation Trust. GM acknowledges support from a UKRI AI Turing Acceleration Fellowship (EPSRC EP/V024868/1).

## References

- Gaon An, Seungyong Moon, Jang-Hyun Kim, and Hyun Oh Song. Uncertainty-based offline reinforcement learning with diversified Q-ensemble. *Advances in neural information processing systems*, 34:7436–7447, 2021. 2, 4
- Arthur Argenson and Gabriel Dulac-Arnold. Model-based offline planning. In *International Conference on Learning Representations*, 2020. 2
- Chenjia Bai, Lingxiao Wang, Zhuoran Yang, Zhi-Hong Deng, Animesh Garg, Peng Liu, and Zhaoran Wang. Pessimistic bootstrapping for uncertainty-driven offline reinforcement learning. In *International Conference on Learning Representations*, 2022. 2, 4
- Alex Beeson and Giovanni Montana. Balancing policy constraint and ensemble size in uncertainty-based offline reinforcement learning. *Machine Learning*, 113(1):443–488, 2024. 2, 4, 14
- David Brandfonbrener, Will Whitney, Rajesh Ranganath, and Joan Bruna. Offline RL without off-policy evaluation. *Advances in neural information processing systems*, 34:4933–4946, 2021. 4, 9
- Yash Chandak, Georgios Theodorou, James Kostas, Scott Jordan, and Philip Thomas. Learning action representations for reinforcement learning. In *International conference on machine learning*, pp. 941–950. PMLR, 2019. 9, 10, 23
- Yevgen Chebotar, Quan Vuong, Karol Hausman, Fei Xia, Yao Lu, Alex Irpan, Aviral Kumar, Tianhe Yu, Alexander Herzog, Karl Pertsch, et al. Q-transformer: Scalable offline reinforcement learning via autoregressive Q-functions. In *Conference on Robot Learning*, pp. 3909–3928. PMLR, 2023. 4
- Caroline Claus and Craig Boutilier. The dynamics of reinforcement learning in cooperative multiagent systems. *AAAI/IAAI*, 1998(746-752):2, 1998. 3
- Driess, Ha, and Toussaint. Deep visual reasoning - learning to predict action sequences for assembly tasks. In *2020 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, pp. 4645–4651. IEEE, 2020. 2
- Wei Du, Shifei Ding, Lili Guo, Jian Zhang, Chenglong Zhang, and Ling Ding. Value function factorization with dynamic weighting for deep multi-agent reinforcement learning. *Information Sciences*, 615:191–208, 2022. 3, 5
- Gabriel Dulac-Arnold, Richard Evans, Hado van Hasselt, Peter Sunehag, Timothy Lillicrap, Jonathan Hunt, Timothy Mann, Theophane Weber, Thomas Degris, and Ben Coppin. Deep reinforcement learning in large discrete action spaces. *arXiv preprint arXiv:1512.07679*, 2015. 4
- Gregory Farquhar, Laura Gustafson, Zeming Lin, Shimon Whiteson, Nicolas Usunier, and Gabriel Synnaeve. Growing action spaces. In *International Conference on Machine Learning*, pp. 3040–3051. PMLR, 2020. 4
- Justin Fu, Aviral Kumar, Ofir Nachum, George Tucker, and Sergey Levine. D4RL: Datasets for deep data-driven reinforcement learning. *arXiv preprint arXiv:2004.07219*, 2020. 4, 9
- Scott Fujimoto and Shixiang Shane Gu. A minimalist approach to offline reinforcement learning. *Advances in Neural Information Processing Systems*, 34:20132–20145, 2021. 2, 4
- Scott Fujimoto, Herke Hoof, and David Meger. Addressing function approximation error in actor-critic methods. In *International conference on machine learning*, pp. 1587–1596. PMLR, 2018. 4
- Scott Fujimoto, Edoardo Conti, Mohammad Ghavamzadeh, and Joelle Pineau. Benchmarking batch deep reinforcement learning algorithms. *arXiv preprint arXiv:1910.01708*, 2019a. 4, 8
- Scott Fujimoto, David Meger, and Doina Precup. Off-policy deep reinforcement learning without exploration. In *International Conference on Machine Learning*, pp. 2052–2062. PMLR, 2019b. 1, 2, 3, 4, 8



- Chris Gaskett. *Q-learning for robot control*. PhD thesis, Australian National University, 2002. 3
- Kamyar Ghasemipour, Shixiang Shane Gu, and Ofir Nachum. Why so pessimistic? estimating uncertainties for offline RL through ensembles, and why their independence matters. *Advances in Neural Information Processing Systems*, 35:18267–18281, 2022. 2, 4
- Omer Gottesman, Fredrik Johansson, Joshua Meier, Jack Dent, Donghun Lee, Srivatsan Srinivasan, Linying Zhang, Yi Ding, David Wihl, Xuefeng Peng, et al. Evaluating reinforcement learning algorithms in observational health settings. *arXiv preprint arXiv:1805.12298*, 2018. 5
- Pengjie Gu, Mengchen Zhao, Chen Chen, Dong Li, Jianye Hao, and Bo An. Learning pseudometric-based action representations for offline reinforcement learning. In *International Conference on Machine Learning*, pp. 7902–7918. PMLR, 2022. 4
- Carlos Guestrin, Daphne Koller, Ronald Parr, and Shobha Venkataraman. Efficient solution algorithms for factored MDPs. *Journal of Artificial Intelligence Research*, 19:399–468, 2003. 7
- Caglar Gulcehre, Sergio Gómez Colmenarejo, Jakub Sygnowski, Thomas Paine, Konrad Zolna, Yutian Chen, Matthew Hoffman, Razvan Pascanu, Nando de Freitas, et al. Addressing extrapolation error in deep offline reinforcement learning. *Offline Reinforcement Learning Workshop at Neural Information Processing Systems, 2020*, 2020a. 4
- Caglar Gulcehre, Ziyu Wang, Alexander Novikov, Thomas Paine, Sergio Gómez, Konrad Zolna, Rishabh Agarwal, Josh S Merel, Daniel J Mankowitz, Cosmin Paduraru, et al. RL unplugged: A suite of benchmarks for offline reinforcement learning. *Advances in Neural Information Processing Systems*, 33:7248–7259, 2020b. 4, 9
- Matteo Hessel, Joseph Modayil, Hado Van Hasselt, Tom Schaul, Georg Ostrovski, Will Dabney, Dan Horgan, Bilal Piot, Mohammad Azar, and David Silver. Rainbow: Combining improvements in deep reinforcement learning. In *Thirty-second AAAI conference on artificial intelligence*, 2018. 1, 3
- Dan Horgan, John Quan, David Budden, Gabriel Barth-Maron, Matteo Hessel, Hado van Hasselt, and David Silver. Distributed prioritized experience replay. In *International Conference on Learning Representations*, 2018. 24
- Thomas Hubert, Julian Schrittwieser, Ioannis Antonoglou, Mohammadamin Barekatain, Simon Schmitt, and David Silver. Learning and planning in complex action spaces. In *International Conference on Machine Learning*, pp. 4476–4486. PMLR, 2021. 4
- David Ireland and Giovanni Montana. Revalued: Regularised ensemble value-decomposition for factorisable Markov decision processes. In *The Twelfth International Conference on Learning Representations*, 2023. 5, 7, 9, 21, 24, 25, 28, 29
- Michael Janner, Yilun Du, Joshua Tenenbaum, and Sergey Levine. Planning with diffusion for flexible behavior synthesis. In *International Conference on Machine Learning*, pp. 9902–9915. PMLR, 2022. 2
- Dmitry Kalashnikov, Alex Irpan, Peter Pastor, Julian Ibarz, Alexander Herzog, Eric Jang, Deirdre Quillen, Ethan Holly, Mrinal Kalakrishnan, Vincent Vanhoucke, et al. Scalable deep reinforcement learning for vision-based robotic manipulation. In *Conference on robot learning*, pp. 651–673. PMLR, 2018. 1
- Michael Kearns and Daphne Koller. Efficient reinforcement learning in factored MDPs. In *IJCAI*, volume 16, pp. 740–747, 1999. 7
- Rahul Kidambi, Aravind Rajeswaran, Praneeth Netrapalli, and Thorsten Joachims. MOREL: Model-based offline reinforcement learning. *Advances in neural information processing systems*, 33:21810–21823, 2020. 2
- Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014. 24, 25

- B Ravi Kiran, Ibrahim Sobh, Victor Talpaert, Patrick Mannion, Ahmad A. Al Sallab, Senthil Yogamani, and Patrick Pérez. Deep reinforcement learning for autonomous driving: A survey. *IEEE Transactions on Intelligent Transportation Systems*, 23(6):4909–4926, 2022. 1
- Ilya Kostrikov, Rob Fergus, Jonathan Tompson, and Ofir Nachum. Offline reinforcement learning with Fisher divergence critic regularization. In *International Conference on Machine Learning*, pp. 5774–5783. PMLR, 2021a. 2, 4
- Ilya Kostrikov, Ashvin Nair, and Sergey Levine. Offline reinforcement learning with implicit Q-Learning. In *International Conference on Learning Representations*, 2021b. 2, 4, 9
- Landon Kraemer and Bikramjit Banerjee. Multi-agent reinforcement learning as a rehearsal for decentralized planning. *Neurocomputing*, 190:82–94, 2016. 5
- Aviral Kumar, Justin Fu, Matthew Soh, George Tucker, and Sergey Levine. Stabilizing off-policy Q-learning via bootstrapping error reduction. *Advances in neural information processing systems*, 32, 2019. 2, 4
- 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. 2, 8
- Sascha Lange, Thomas Gabel, and Martin Riedmiller. Batch reinforcement learning. *Springer*, pp. 45–73, 2012. 1, 3
- Sergey Levine, Aviral Kumar, George Tucker, and Justin Fu. Offline reinforcement learning: Tutorial, review, and perspectives on open problems. *arXiv preprint arXiv:2005.01643*, 2020. 1, 2, 3
- Kaixiang Lin, Renyu Zhao, Zhe Xu, and Jiayu Zhou. Efficient large-scale fleet management via multi-agent deep reinforcement learning. In *Proceedings of the 24th ACM SIGKDD international conference on knowledge discovery & data mining*, pp. 1774–1783, 2018. 2
- Siqi Liu, Kay Choong See, Kee Yuan Ngiam, Leo Anthony Celi, Xingzhi Sun, and Mengling Feng. Reinforcement learning for clinical decision support in critical care: comprehensive review. *Journal of medical Internet research*, 22(7):e18477, 2020. 2
- Jianlan Luo, Perry Dong, Jeffrey Wu, Aviral Kumar, Xinyang Geng, and Sergey Levine. Action-quantized offline reinforcement learning for robotic skill learning. In *7th Annual Conference on Robot Learning*, 2023. 8, 9
- A Rupam Mahmood, Dmytro Korenkevych, Gautham Vasan, William Ma, and James Bergstra. Benchmarking reinforcement learning algorithms on real-world robots. In *Conference on robot learning*, pp. 561–591. PMLR, 2018. 1
- Luke Metz, Julian Ibarz, Navdeep Jaitly, and James Davidson. Discrete sequential prediction of continuous actions for deep RL. *arXiv preprint arXiv:1705.05035*, 2017. 4
- Volodymyr Mnih, Koray Kavukcuoglu, David Silver, Alex Graves, Ioannis Antonoglou, Daan Wierstra, and Martin Riedmiller. Playing atari with deep reinforcement learning. *arXiv preprint arXiv:1312.5602*, 2013. 1, 3
- Alexander Nikulin, Vladislav Kurenkov, Denis Tarasov, and Sergey Kolesnikov. Anti-exploration by random network distillation. In *International Conference on Machine Learning*, pp. 26228–26244. PMLR, 2023. 2
- Thomas Pierrot, Valentin Macé, Jean-Baptiste Sevestre, Louis Monier, Alexandre Laterre, Nicolas Perrin, Karim Beguir, and Olivier Sigaud. Factored action spaces in deep reinforcement learning. 2021. 4
- Tabish Rashid, Gregory Farquhar, Bei Peng, and Shimon Whiteson. Weighted qmix: Expanding monotonic value function factorisation for deep multi-agent reinforcement learning. *Advances in neural information processing systems*, 33:10199–10210, 2020a. 3, 5

- Tabish Rashid, Mikayel Samvelyan, Christian Schroeder De Witt, Gregory Farquhar, Jakob Foerster, and Shimon Whiteson. Monotonic value function factorisation for deep multi-agent reinforcement learning. *Journal of Machine Learning Research*, 21(178):1–51, 2020b. 3, 5
- Aaman Peter Rebello, Shengpu Tang, Jenna Wiens, and Sonali Parbhoo. Leveraging factored action spaces for off-policy evaluation. In *ICML Workshop on New Frontiers in Learning, Control, and Dynamical Systems*, 2023. 14
- Tim Seyde, Igor Gilitschenski, Wilko Schwarting, Bartolomeo Stellato, Martin Riedmiller, Markus Wulfmeier, and Daniela Rus. Is bang-bang control all you need? solving continuous control with Bernoulli policies. *Advances in Neural Information Processing Systems*, 34:27209–27221, 2021. 4
- Tim Seyde, Peter Werner, Wilko Schwarting, Igor Gilitschenski, Martin Riedmiller, Daniela Rus, and Markus Wulfmeier. Solving continuous control via Q-learning. In *The Eleventh International Conference on Learning Representations*, 2022. 2, 3, 5, 9, 24, 25, 29
- Sahil Sharma, Aravind Suresh, Rahul Ramesh, and Balaraman Ravindran. Learning to factor policies and action-value functions: Factored action space representations for deep reinforcement learning. *arXiv preprint arXiv:1705.07269*, 2017. 4
- Peter Sunehag, Guy Lever, Audrunas Gruslys, Wojciech Marian Czarnecki, Vinicius Zambaldi, Max Jaderberg, Marc Lanctot, Nicolas Sonnerat, Joel Z Leibo, Karl Tuyls, et al. Value-decomposition networks for cooperative multi-agent learning. *arXiv preprint arXiv:1706.05296*, 2017. 3, 5
- Richard S Sutton and Andrew G Barto. Reinforcement learning: An introduction. *MIT press*, 2018. 2, 3
- Phillip Swazinna, Steffen Udluft, and Thomas Runkler. Overcoming model bias for robust offline deep reinforcement learning. *Engineering Applications of Artificial Intelligence*, 104:104366, 2021. 2
- Shengpu Tang, Maggie Makar, Michael Sjoding, Finale Doshi-Velez, and Jenna Wiens. Leveraging factored action spaces for efficient offline reinforcement learning in healthcare. *Advances in Neural Information Processing Systems*, 35:34272–34286, 2022. 5
- Yunhao Tang and Shipra Agrawal. Discretizing continuous action space for on-policy optimization. In *Proceedings of the aaai conference on artificial intelligence*, volume 34, pp. 5981–5988, 2020. 2, 4
- Arash Tavakoli, Fabio Pardo, and Petar Kormushev. Action branching architectures for deep reinforcement learning. In *Proceedings of the aaai conference on artificial intelligence*, volume 32, 2018. 2, 3, 4, 28
- Sebastian Thrun and Anton Schwartz. Issues in using function approximation for reinforcement learning. In *Proceedings of the Fourth Connectionist Models Summer School*, volume 255, pp. 263. Hillsdale, NJ, 1993. 3, 5, 21
- Saran Tunyasuvunakool, Alistair Muldal, Yotam Doron, Siqu Liu, Steven Bohez, Josh Merel, Tom Erez, Timothy Lillicrap, Nicolas Heess, and Yuval Tassa. dm\_control: Software and tasks for continuous control. *Software Impacts*, 6:100022, 2020. 9
- Tom Van de Wiele, David Warde-Farley, Andriy Mnih, and Volodymyr Mnih. Q-learning in enormous action spaces via amortized approximate maximization. *arXiv preprint arXiv:2001.08116*, 2020. 4
- Yifan Wu, George Tucker, and Ofir Nachum. Behavior regularized offline reinforcement learning. *arXiv preprint arXiv:1911.11361*, 2019. 2, 4
- Rui Yang, Chenjia Bai, Xiaoteng Ma, Zhaoran Wang, Chongjie Zhang, and Lei Han. RORL: Robust offline reinforcement learning via conservative smoothing. In *Advances in Neural Information Processing Systems*, 2022. 2, 4
- Chao Yu, Jiming Liu, Shamim Nemati, and Guosheng Yin. Reinforcement learning in healthcare: A survey. *ACM Computing Surveys (CSUR)*, 55(1):1–36, 2021a. 1

Tianhe Yu, Garrett Thomas, Lantao Yu, Stefano Ermon, James Y Zou, Sergey Levine, Chelsea Finn, and Tengyu Ma. MOPO: Model-based offline policy optimization. *Advances in Neural Information Processing Systems*, 33:14129–14142, 2020. 2

Tianhe Yu, Aviral Kumar, Rafael Rafailov, Aravind Rajeswaran, Sergey Levine, and Chelsea Finn. COMBO: Conservative offline model-based policy optimization. *Advances in neural information processing systems*, 34:28954–28967, 2021b. 2

Xiangyu Zhao, Long Xia, Liang Zhang, Zhuoye Ding, Dawei Yin, and Jiliang Tang. Deep reinforcement learning for page-wise recommendations. In *Proceedings of the 12th ACM conference on recommender systems*, pp. 95–103, 2018. 2

Wenxuan Zhou, Sujay Bajracharya, and David Held. PLAS: Latent action space for offline reinforcement learning. In *Conference on Robot Learning*, pp. 1719–1735. PMLR, 2021. 2

## Appendix

### A DQN vs DecQN for uniformly distributed errors

To provide support for the ideas presented in Section 4, we investigate properties of the target difference under DQN and DecQN based on uniformly distribution noise.

As demonstrated by [Thrun & Schwartz \(1993\)](#) and [Ireland & Montana \(2023\)](#), if  $\epsilon(s, \mathbf{a})$  are modelled as independent identically distributed (i.i.d.) uniform random variables  $U(-b, b)$ , then for  $|A|$  actions the expectation and variance of the target difference are, respectively,

$$\mathbb{E}[Z_s] = \gamma b \frac{|A| - 1}{|A| + 1}, \quad (9)$$

$$\text{Var}(Z_s) = \gamma 4b^2 \frac{|A|}{(|A| + 1)^2(|A| + 2)}. \quad (10)$$

For  $|A| > 1$ , i.e. more than one action,  $\mathbb{E}[Z_s]$  is positive and increasing w.r.t.  $|A|$  whereas  $\text{Var}(Z_s)$  is positive and decreasing. Hence there is a bias/variance trade-off for overestimation in target Q-values.

This trade-off becomes apparent when comparing overestimation for DQN and DecQN, where  $|A| = \prod_{i=1}^N n_i$  for the former and  $|A| = \sum_{i=1}^N n_i$  for the latter. As shown by [Ireland & Montana \(2023\)](#), if  $\epsilon^i(s, a_i)$  are also modelled as i.i.d. uniform random variables  $U(-b, b)$ , then the expected target difference under DecQN is lower than DQN ( $\mathbb{E}[Z_{s'}^{dec}] \leq \mathbb{E}[Z_{s'}^{dqn}]$ ), but the variance is higher ( $V[Z_{s'}^{dec}] \geq V[Z_{s'}^{dqn}]$ ).

We now consider this bias/variance trade-off in the context of in-distribution and out-of-distribution actions/sub-actions. This time, let  $U(-b, b)$  and  $U(-kb, kb)$  be the distribution of errors for in-distribution and out-of-distribution actions/sub-actions, respectively, where  $k > 1$ . When all actions/sub-actions are either in-distribution or out-of-distribution the above conclusions hold, with the additional observation that both bias and variance are lower when all actions are in-distribution compared to all out-of-distribution (as  $b < kb$ ).

The picture becomes more complex when there is a mixture of in-distribution and out-of-distribution actions/sub-actions. This stems from the fact there is no longer a closed form for the expectation/variance of the target difference since errors arise from two different distributions. This is exacerbated by the relative coverage of sub-actions and atomic actions. In terms of expectation, since this is an increasing function of both  $b$  and  $|A|$  we can deduce the property  $\mathbb{E}[Z_{s'}^{dec}] \leq \mathbb{E}[Z_{s'}^{dqn}]$  holds regardless of the value of  $k$  or the number of in-distribution actions/sub-actions. The variance is more nuanced since this is an increasing function for  $b$  but decreasing for  $|A|$ . The value of  $k$  and number of in-distribution actions/sub-actions will determine whether the property  $V[Z_{s'}^{dec}] \geq V[Z_{s'}^{dqn}]$  still holds.

To illustrate this, we conduct a series of simulations based on different action/sub-action configurations and compare how the expectation and variance of the target difference under DQN and DecQN change based on the coverage of atomic actions. The details are as follows.

**DQN simulations:** Let  $|A^{in}|$  and  $|A^{out}|$  be the number of in-distribution and out-of-distribution atomic actions for a given state, respectively, such that  $|A| = |A^{in}| + |A^{out}|$ . For  $|A^{in}| \in \{0, 1, \dots, |A|\}$  sample  $|A^{in}|$  values from a  $U(-b, b)$  distribution and  $|A| - |A^{in}|$  values from a  $U(-kb, kb)$  distribution and store the maximum value for the pooled sample. Repeat this 10000 times and calculate the mean and variance of these stored maximums. These are the estimates of  $\mathbb{E}[Z_s^{dqn}]$  and  $V[Z_s^{dqn}]$ .

**DecQN simulations** Let  $|A_i^{in}|$  and  $|A_i^{out}|$  be the number of in-distribution and out-of-distribution sub-actions in dimension  $i$  for a given state, respectively, such that  $|A_i| = |A_i^{in}| + |A_i^{out}|$ . For  $|A^{in}| \in \{0, 1, \dots, |A|\}$  sample  $|A^{in}|$  atomic action from all available actions without replacement and calculate the number of factorised actions  $|A_i^{in}|$  in each sub-action dimension  $i$ . In each sub-action dimension  $i$  sample  $|A_i^{in}|$  values from a  $U(-b, b)$  distribution and  $|A_i| - |A_i^{in}|$  values from a  $U(-kb, kb)$  distribution and store the maximum value for the pooled sample for each dimension  $i$ . Calculate and record the average maximum value across all dimensions  $N$  (as per the DecQN decomposition). Repeat this 10000 times and calculate the mean and

variance for these stored maximums. Since sub-action coverage is dependent on which atomic actions are sampled, we repeat this entire procedure 100 times to sample a broad range of atomic actions and calculate the overall mean and variance across these 100 simulations. These are the estimates of  $\mathbb{E}[Z_s^{dec}]$  and  $V[Z_s^{dec}]$ .

In Figure 8 we summarise these simulations for three configurations:  $N = 3$  each with two sub-actions,  $N = 4$  each with two sub-actions and  $N = 3$  each with three sub-actions. We set  $b = 1$  and  $k = 2$ . For the expectation, we see that in all cases  $\mathbb{E}[Z_{s'}^{dec}] \leq \mathbb{E}[Z_{s'}^{dqn}]$  for all values of  $|A^{in}|$ . For the variance, we see that in all cases  $V[Z_{s'}^{dec}] \geq V[Z_{s'}^{dqn}]$  for values of  $|A^{in}|$  close to 0 and  $|A|$  but in between  $V[Z_{s'}^{dec}] \leq V[Z_{s'}^{dqn}]$ . This aligns with our intuition, namely that the expectation of the target difference under DecQN will always be lower than DQN, but the variance will depend on the relative proportions of in-distribution actions/sub-actions under atomic and factorised representations. Note that if all atomic actions are in-distribution or out-of-distribution (i.e.  $|A^{in}| = 0$  or  $|A^{in}| = |A|$ ) the expectation and variance revert back to Equations 9 and 10.

## B Algorithms - full procedures

This section contains the full procedures for algorithms detailed in Section 5.

---

### Algorithm 1 DecQN-BCQ

**Require:** Threshold  $\tau$ , discount factor  $\gamma$ , target network update rate  $\mu$ , number sub-action spaces  $N$  and dataset  $\mathcal{B}$ .

Initialise utility function parameters  $\theta = \{\theta_i\}_{i=1}^N$ , corresponding target parameters  $\hat{\theta} = \theta$  and policy parameters  $\phi = \{\phi_i\}_{i=1}^N$

**for**  $t = 0$  to  $T$  **do**

    Sample minibatch of transitions  $(s, \mathbf{a}, r, s')$  from  $\mathcal{B}$

$\phi \leftarrow \arg \min_{\phi} \frac{1}{N} \sum_{i=1}^N - \sum_{s, a_i} \log \pi_{\phi_i}^i(a_i | s)$

$\theta \leftarrow \arg \min_{\theta} \sum_{s, \mathbf{a}, r, s'} (Q_{\theta}(s, \mathbf{a}) - y)^2$

    where:

$Q_{\theta}(s, \mathbf{a}) = 1/N \sum_{i=1}^N U_{\theta_i}^i(s, a_i),$

$y = r + \gamma/N \sum_{i=1}^N \max_{a'_i : \rho^i(a'_i) \geq \tau} U_{\theta_i}^i(s', a'_i),$

$\rho^i(a'_i) = \pi_{\phi_i}^i(a'_i | s') / \max_{\hat{a}'_i} \pi_{\phi_i}^i(\hat{a}'_i | s')$

$\hat{\theta} \leftarrow \mu\theta + (1 - \mu)\hat{\theta}$

**end for**

---

### Algorithm 2 DecQN-CQL

**Require:** Conservative coefficient  $\alpha$ , discount factor  $\gamma$ , target network update rate  $\mu$ , number sub-action spaces  $N$  and dataset  $\mathcal{B}$ .

Initialise utility function parameters  $\theta = \{\theta_i\}_{i=1}^N$  and corresponding target parameters  $\hat{\theta} = \theta$

**for**  $t = 0$  to  $T$  **do**

    Sample minibatch of transitions  $(s, \mathbf{a}, r, s')$  from  $\mathcal{B}$

$\theta \leftarrow \arg \min_{\theta} \sum_{s, \mathbf{a}, r, s'} (Q_{\theta}(s, \mathbf{a}) - y)^2 + \alpha \sum_{s, \mathbf{a}} \frac{1}{N} \sum_{i=1}^N [\log \sum_{a_j \in \mathcal{A}_i} \exp(U_{\theta_i}^i(s, a_j)) - U_{\theta_i}^i(s, a_i)]$

    where:

$Q_{\theta}(s, \mathbf{a}) = 1/N \sum_{i=1}^N U_{\theta_i}^i(s, a_i),$

$y = r + 1/N \sum_{i=1}^N \max_{a'_i} U_{\theta_i}^i(s', a'_i)$

$\hat{\theta} \leftarrow \mu\theta + (1 - \mu)\hat{\theta}$

**end for**

---

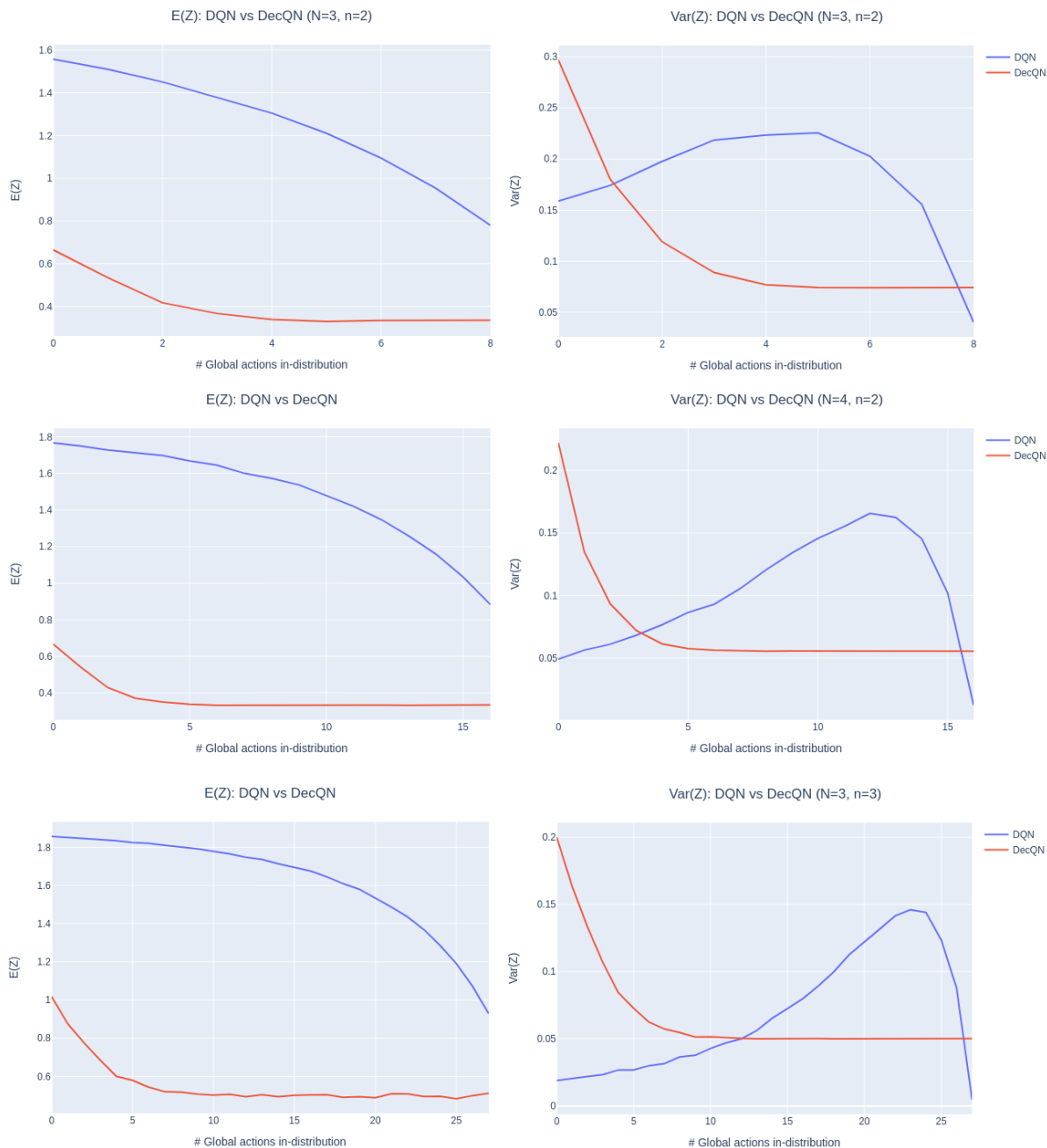


Figure 8: Comparing the expectation and variance of target differences under DQN and DecQN for various configurations and coverages of atomic/sub-actions.

## C Data collection procedures

### C.1 Maze

For the Maze environment, we used the default settings as per (Chandak et al., 2019) with the exception that all actions are available from the beginning. To collect the datasets, we first trained agents for varying numbers of actuators using DecQN, parameterising utility functions as 2-layer MLPs each with 512 nodes

**Algorithm 3** DecQN-IQL

**Require:** Expectile  $\tau$ , discount factor  $\gamma$ , target network update rate  $\mu$ , number sub-action spaces  $N$  and dataset  $\mathcal{B}$ .

Initialise utility function parameters  $\theta = \{\theta_i\}_{i=1}^N$  and corresponding target parameters  $\hat{\theta} = \theta$ . Initialise state value function parameters  $\psi$  and policy parameters  $\phi = \{\phi_i\}_{i=1}^N$

**for**  $t = 0$  to  $T$  **do**

Sample minibatch of transitions  $(s, \mathbf{a}, r, s')$  from  $\mathcal{B}$

$$\phi \leftarrow \arg \min_{\phi} \frac{1}{N} \sum_{i=1}^N - \sum_{s, a_i} \log \pi_{\phi_i}^i(a_i | s)$$

$$\theta \leftarrow \arg \min_{\theta} \sum_{s, \mathbf{a}, r, s'} (Q_{\theta}(s, \mathbf{a}) - y)^2$$

$$\psi \leftarrow \arg \min_{\psi} \sum_{s, \mathbf{a}} [L_2^{\tau}(Q_{\hat{\theta}}(s, \mathbf{a}) - V_{\psi}(s))]$$

where:

$$Q_{\theta}(s, \mathbf{a}) = 1/N \sum_{i=1}^N U_{\theta_i}^i(s, a_i),$$

$$Q_{\hat{\theta}}(s, \mathbf{a}) = 1/N \sum_{i=1}^N U_{\hat{\theta}_i}^i(s, a_i),$$

$$y = r + V_{\psi}(s')$$

$$\hat{\theta} \leftarrow \mu\theta + (1 - \mu)\hat{\theta}$$

**end for**

**Algorithm 4** DecQN-OneStep

**Require:** Discount factor  $\gamma$ , target network update rate  $\mu$ , number sub-action spaces  $N$  and dataset  $\mathcal{B}$ .

Initialise utility function parameters  $\theta = \{\theta_i\}_{i=1}^N$  and corresponding target parameters  $\hat{\theta} = \theta$ . Initialise policy parameters  $\phi = \{\phi_i\}_{i=1}^N$

**for**  $t = 0$  to  $T$  **do**

Sample minibatch of transitions  $(s, \mathbf{a}, r, s')$  from  $\mathcal{B}$

$$\phi \leftarrow \arg \min_{\phi} \frac{1}{N} \sum_{i=1}^N - \sum_{s, a_i} \log \pi_{\phi_i}^i(a_i | s)$$

$$\theta \leftarrow \arg \min_{\theta} \sum_{s, \mathbf{a}, r, s'} (Q_{\theta}(s, \mathbf{a}) - y)^2$$

where:

$$Q_{\theta}(s, \mathbf{a}) = 1/N \sum_{i=1}^N U_{\theta_i}^i(s, a_i),$$

$$y = r + 1/N \sum_{i=1}^N \sum_{a_i} \pi_{\phi_i}^i(a_i | s) U_{\hat{\theta}_i}^i(s', a'_i)$$

$$\hat{\theta} \leftarrow \mu\theta + (1 - \mu)\hat{\theta}$$

**end for**

that take in a state and output sub-action utility values. We maintained a fixed exploration parameter  $\epsilon = 0.1$  and updated target network parameters using Polyak-averaging with  $\mu = 0.005$ . We used the Adam optimiser (Kingma & Ba, 2014) with learning rate  $3e - 4$  and a batch size of 256.

Once we have trained the DecQNs we create the benchmark datasets by collecting data using a greedy policy derived from the learned utility values. Each dataset contains 10k transitions. For the expert policy we trained the DecQNs until their test performance approached the maximum return possible (approximately 100). For the medium policy we aimed for a test performance of approximately 1/3rd of the expert.

**C.2 DeepMind Control Suite**

To collect the DMC suite datasets we followed the training procedures laid out by (Seyde et al., 2022; Ireland & Montana, 2023) to train the Decoupled Q-networks. To expedite the data collection process during training we used a distributed setup using multiple workers to collect data in parallel (Horgan et al., 2018). We parameterised the utility functions using a (shared) single ResNet layer followed by layer norm, followed by a linear head for each of the sub-action spaces which predicts sub-action utility values. Full details regarding network architecture and hyperparameters can be found in Table 1.

Once we have trained the DecQNs we create the benchmark datasets by collecting data using a greedy policy derived from the learned utility values. For the case study in Section 7.1 each dataset contains 100k transitions and for the main experiments in Section 7.3 each dataset contains 1M transitions. As each episode



is truncated at 1,000 time steps in the DM control suite, this corresponds to collecting 100 episodes for the case study and 1,000 episodes for the main experiments. For the expert policy we trained the DecQNs until their test performance corresponded to the performance given in (Seyde et al., 2022; Ireland & Montana, 2023). For the medium policy we aimed for a test performance of approximately 1/3rd of the reported expert score.

We largely employ the same hyperparameters as the original DecQN study, as detailed in Table 1. Exceptions include the decay of the exploration parameter ( $\epsilon$ ) to a minimum value instead of keeping it constant, and the use of Polyak-averaging for updating the target network parameters, as opposed to a hard reset after every specified number of updates. Finally, we sample from the replay buffer uniformly at random, as opposed to using a priority. We maintain the same hyperparameters across all our experiments.

For  $n = 3$  we use DecQN to train networks and collect datasets. For  $n > 3$  we use REValueD to train networks and collect datasets due to better scaling to higher numbers of bins (Ireland & Montana, 2023).

Table 1: Hyperparameters used in DecQN and REValueD training.

Parameters	Value
Optimizer	Adam
Learning rate	$1 \times 10^{-4}$
Replay size	$5 \times 10^5$
n-step returns	3
Discount, $\gamma$	0.99
Batch size	256
Hidden size	512
Gradient clipping	40
Target network update parameter, $c$	0.005
Imp. sampling exponent	0.2
Priority exponent	0.6
Minimum exploration, $\epsilon$	0.05
$\epsilon$ decay rate	0.99995
Regularisation loss coefficient $\beta$	0.5
Ensemble size K	10

## D Full implementation details

For both Maze and DMC environments/tasks, utility functions are parameterised by neural networks, comprising a 2-layer MLP with ReLU activation functions and 512 nodes, taking in a normalised state as input and outputting utility values for each sub-action space. We use the same architecture for policies, with the output layer a softmax across actions within each sub-actions-space. State value functions mirror this architecture except in the final layer which outputs a single value. We train networks via stochastic gradient descent using the Adam optimiser (Kingma & Ba, 2014) with learning rate  $3e^{-4}$  and a batch size of 256. For state and state-action value functions we use the Huber loss as opposed to MSE loss. We set the discount factor  $\gamma = 0.99$  and the target network update rate  $\mu = 0.005$ . We utilise a dual-critic approach, taking the mean across two utility estimates for target Q-values. For the maze task, agents are trained for 100k gradient updates. For the DMC tasks, agents are trained for 1M gradient updates.

The only hyperparameters we tune are the threshold  $\tau$  in BCQ, conservative coefficient  $\alpha$  in CQL, expectile  $\tau$  and balance coefficient  $\lambda$  in IQL and balance coefficient  $\lambda$  in OneStep. We allow these to vary across environment/task, but to better reflect real-world scenarios where the quality of data may be unknown, we forbid variation within environments/tasks.

Table 2 provides hyperparameters for all environments/tasks. For DecQN-BCQ we searched over  $\tau = \{0.025, 0.05, 0.1, 0.25, 0.5, 0.75\}$ . For DecQN-CQL we searched over  $\alpha = \{0.25, 0.5, 1, 2\}$ . For DecQN-IQL

we searched over  $\tau = \{0.5, 0.6, 0.7, 0.8\}$ ,  $\lambda = \{1, 2, 5, 10, 20, 50\}$ . For DecQN-OneStep we searched over  $\lambda = \{1, 2, 5, 10, 20, 50\}$ .

Table 2: Hyperparameters for experiments in Section 7

Environment/task	Number of bins ( $n$ )	BCQ $\tau$	CQL $\alpha$	IQL $\beta, \lambda$	OneStep $\lambda$
Maze	3	0.5	1	0.5, 20	50
Maze	5	0.5	0.25	0.5, 20	20
Maze	7	0.5	0.5	0.5, 50	20
Maze	10	0.5	0.5	0.5, 50	20
Maze	12	0.5	0.5	0.5, 50	50
Maze	15	0.5	0.5	0.5, 20	20
FingerSpin	3	0.5	0.25	0.5, 1	1
FishSwim	3	0.25	0.25	0.5, 1	1
CheetahRun	3	0.05	0.25	0.5, 1	1
QuadrupedWalk	3	0.25	0.25	0.5, 2	1
HumanoidStand	3	0.5	0.25	0.5, 2	2
DogTrot	3	0.25	1	0.5, 5	5
DogTrot	10	0.25	0.25	0.5, 5	2
DogTrot	30	0.75	1	0.5, 2	2
DogTrot	50	0.5	0.5	0.5, 2	2
DogTrot	75	0.7	0.5	0.5, 2	2
DogTrot	100	0.5	2	0.5, 5	5

### D.1 Allowing hyperparameter variation within environment/task

As per Section 8, in Table 3 we provide examples of performance improvement after permitting hyperparameter variation within the same environment/task. In general, we see lower quality datasets benefit from smaller hyperparameters (i.e. those that weight more towards RL and less towards BC) and higher quality datasets benefit from larger hyperparameters (i.e. those the weight more towards BC and less towards RL). This mirrors findings from previous papers outlined in Section 3.

## E Number of actions under atomic and factorised representations

Table 4 summarises the number of actions requiring value estimation under atomic and factorised representations.

## F Case study Q-value errors

To provide further insights into Q-value errors for the case study in Section 7.1, we examine the evolution of Q-value errors during training. Every 5k gradient updates we obtain a MC estimate of true Q-values using discounted rewards from environmental rollouts. We then compare these MC estimates with Q-values predicted by both DQN-CQL and DecQN-CQL networks for the respective actions taken. To make better use of rollouts (which can be expensive), we calculate MC estimates and DQN-/DecQN-CQL Q-values for the first 500 states in the trajectory, as using a discount factor of  $\gamma = 0.99$  discounts rewards by over 99% for time-steps beyond 500 (and all tasks considered have trajectory length 1000). In total we perform 10 rollouts, giving 5000 estimates of the error between true and DQN-CQL/DecQN-CQL Q-values. In Figure 9 we plot the mean absolute error over the course of training, with the solid line representing the mean across five random seeds and shaded area the standard error. For all values of  $n_i$  we observe the mean absolute error is less for DecQN-CQL than DQN-CQL, particularly for  $n_i > 3$ , aligning with each algorithm’s respective performance in Figure 3.

Table 3: Individual performance allowing hyperparameters to vary within environment/task; ‘dog-trot’,  $n = 3$ . Figures are mean normalised scores, with 0 and 100 representing random and expert policies, respectively. Highest score highlighted in bold

Environment -dataset				
DogTrot (BCQ)	$\tau = 0.025$	$\tau = 0.05$	$\tau = 0.1$	$\tau = 0.25$
-expert	57.8	85	93.6	<b>94.4</b>
-medium-expert	3.9	10.1	42.7	<b>74.9</b>
-medium	39.8	38.8	34.2	<b>49.1</b>
-random-medium-expert	5	5.6	<b>9</b>	0.1
DogTrot (CQL)	$\alpha = 0.25$	$\alpha = 0.5$	$\alpha = 1$	$\alpha = 2$
-expert	90.8	95.7	99.5	<b>100.2</b>
-medium-expert	76.6	81.7	<b>84.8</b>	75.1
-medium	<b>50.6</b>	48.3	46.5	45.2
-random-medium-expert	41.2	40.8	<b>43.4</b>	38.6
DogTrot (IQL $\tau = 0.5$ )	$\lambda = 1$	$\lambda = 2$	$\lambda = 5$	$\lambda = 10$
-expert	37.9	82.5	98.9	<b>99.5</b>
-medium-expert	33	64	89.3	<b>98.6</b>
-medium	<b>58.8</b>	56.5	52	47.3
-random-medium-expert	10.6	28.6	44.1	<b>44.7</b>
DogTrot (OneStep)	$\lambda = 1$	$\lambda = 2$	$\lambda = 5$	$\lambda = 10$
-expert	53.3	91.7	101.2	<b>102</b>
-medium-expert	44.5	79.4	93.9	<b>96.6</b>
-medium	<b>59.3</b>	57.5	50.2	48
-random-medium-expert	23.8	43.9	44.9	<b>45.1</b>

Table 4: Environment details for DeepMind Control Suite.  $|S|$  represents the size of the state space and  $N$  the number of sub-action spaces.  $\prod_i n_i$  is the total number of actions under atomic representation and  $\sum_i n_i$  under factorised representation when  $n_i = 3$ .

Environment	$ S $	$N$	$\prod_i n_i$	$\sum_i n_i$
Finger Spin	9	2	9	6
Fish Swim	24	5	243	15
Cheetah Run	17	6	729	18
Quadruped Walk	78	12	$\approx 530k$	36
Humanoid Stand	67	21	$\approx 10^{10}$	63
Dog Trot	223	38	$\approx 10^{18}$	114

## G Tabulated results

Tabulated results for Figure 3 are presented in Table 5. Tabulated results for Figure 4 are presented in Table 6. Tabulated results for Figure 5 are presented in Table 7. Tabulated results for Figure 6 are presented in Table 8. Tabulated results for Figure 7 are presented in Table 9.

### Number of seeds

To check five seeds is sufficient for evaluation purposes we run one task for ten seeds and compare performance in Tables 10 and 11, respectively. We see consistent results across different seed counts. We expect similar consistency across all tasks as the same underlying methodology is applied throughout our experiments.

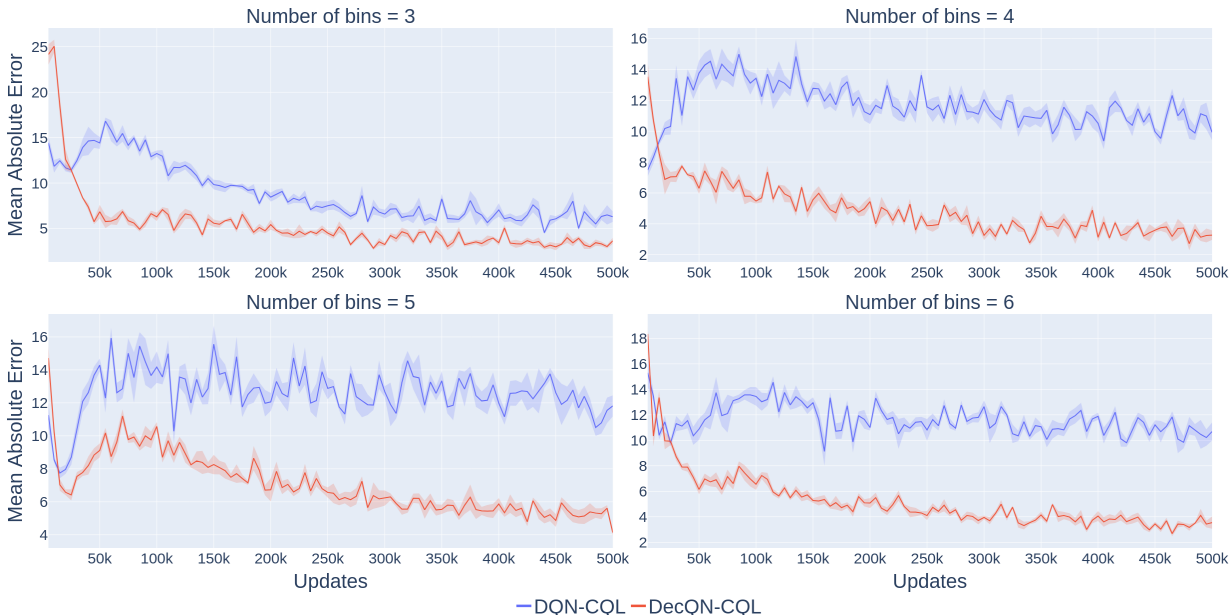


Figure 9: Comparison of estimated errors in Q-values for “cheetah-run-medium-expert” dataset for varying numbers of bins. Errors are lower for DecQN-CQL for all numbers of bins, most notably for  $n_i > 3$ , mirroring the deviation in performance levels between the two approaches.

Table 5: DQN-CQL vs DecQN-CQL - performance and computation comparison for cheetah-run task. Performance figures are mean normalised scores  $\pm$  one standard error, with 0 and 100 representing random and expert policies, respectively. Computation figures are training time and GPU usage. Actions figures are total number of actions requiring value estimation based on atomic/factorised representation.

Method	$n_i$	Actions	Score	Training time (mins)	GPU usage (MB)
DQN-CQL	3	729	79.7 $\pm$ 6.0	19	266
DQN-CQL	4	4096	30.2 $\pm$ 1.7	30	412
DQN-CQL	5	15625	32.2 $\pm$ 1.4	85	958
DQN-CQL	6	46656	32.3 $\pm$ 4.0	227	2388
DecQN-CQL	3	18	92.7 $\pm$ 3.7	19	244
DecQN-CQL	4	24	85.4 $\pm$ 3.8	19	246
DecQN-CQL	5	30	79.5 $\pm$ 1.8	20	246
DecQN-CQL	6	36	84.3 $\pm$ 2.5	20	246

## H Decomposition comparisons

In this Section we compare the DecQN decomposition to two alternative methods that can be used for factorisable discrete action spaces. The first is based on the Branching Dueling Q-Network (BDQ) proposed by Tavakoli et al. (2018). Using our notation, each utility function is considered its own independent Q-function, i.e.

$$Q_{\theta_i}^i(s, a_i) = U_{\theta_i}^i(s, a_i) . \tag{11}$$

Each Q-function is trained by bootstrapping from its own target, and no decomposition is used. That is, the target for  $Q_{\theta_i}^i(s, a_i)$  is given by  $y = r + \gamma \max_{a'_i \in \mathcal{A}_i} Q_{\theta_i}^i(s', a'_i)$ . The findings of Ireland & Montana (2023) demonstrate that, in the online setting, BDQ is unable to match the performance of DecQN. This is likely caused by the fact that, as each sub-action space is now learnt independently, the effects of other sub-actions

Table 6: DQN-CQL vs DecQN-CQL - performance and computation comparison for Maze task with  $N = 15$  actuators. Performance figures are mean normalised scores  $\pm$  one standard error, with 0 and 100 representing random and expert policies, respectively. Computation figures are training time and GPU usage. Actions figures are total number of actions requiring value estimation based on atomic/factorised representation.

Method	Number of transitions	Actions	Score	Training time (mins)	GPU usage (MB)
DQN-CQL	100	32768	0.0 $\pm$ 0.0	34	1728
DQN-CQL	250	32768	48.8 $\pm$ 4.9	34	1728
DQN-CQL	500	32768	68.5 $\pm$ 2.7	34	1728
DQN-CQL	1000	32768	73.7 $\pm$ 2.7	34	1728
DQN-CQL	2500	32768	90.7 $\pm$ 2.5	34	1728
DQN-CQL	5000	32768	96.0 $\pm$ 1.4	34	1728
DecQN-CQL	100	30	31.9 $\pm$ 10.6	4	246
DecQN-CQL	250	30	69.2 $\pm$ 3.2	4	246
DecQN-CQL	500	30	84.9 $\pm$ 2.3	4	246
DecQN-CQL	1000	30	84.8 $\pm$ 1.5	4	246
DecQN-CQL	2500	30	96.4 $\pm$ 1.7	4	246
DecQN-CQL	5000	30	99.7 $\pm$ 0.2	4	246

are treated as effects of the environment dynamics. Due to the fact that each agent is continually updating its own policy, this leads to non-stationary environment dynamics, making the learning problem much more challenging.

We also consider an alternative value-decomposition technique to the mean, namely the sum. That is, we replace the mean operator in Equation 2 with the sum operator:

$$Q_{\theta}(s, \mathbf{a}) = \sum_{i=1}^N U_{\theta_i}^i(s, a_i) . \quad (12)$$

Whilst this may seem a subtle change, Ireland & Montana (2023) proved that the mean and variance of the learning target under this decomposition are both higher than DecQN. Empirical experiments by Seyde et al. (2022); Ireland & Montana (2023) also confirm the inferior performance of the sum decomposition compared to the mean.

In Table 12 we can see that the sum decomposition is less performant in each of the tasks and datasets than the mean. For BDQ, we see that whilst in some cases performance is better than using the sum decomposition, it is generally still less performant than using the mean decomposition. Owing to these results, we focus on the mean decomposition in our main work.

Table 7: Individual performance comparison for Maze for varying numbers of actuators. Figures are mean normalised scores  $\pm$  one standard error, with 0 and 100 representing random and expert policies, respectively. At the bottom of the table we also provide totals, split by data quality and the entire set. We see that DecQN-BCQ is the least performant of the offline methods and that DecQN-CQL/IQL/OneStep perform similarly overall and on expert, medium-expert and medium datasets and DecQN-CQL the best on random-medium-expert datasets.

Environment -dataset	BC	DecQN	DecQN-BCQ	DecQN-CQL	DecQN-IQL	DecQN-OneStep
Maze ( <i>Actuators</i> = 3)						
-expert	100.1 $\pm$ 0	-5.6 $\pm$ 0	100.1 $\pm$ 0	100.2 $\pm$ 0	100.1 $\pm$ 0	100.1 $\pm$ 0
-medium-expert	98.7 $\pm$ 1	-5.6 $\pm$ 0	99.6 $\pm$ 0.6	100.1 $\pm$ 0	100.1 $\pm$ 0	99.9 $\pm$ 0
-medium	55.3 $\pm$ 2.5	-5.6 $\pm$ 0	54.7 $\pm$ 1.7	62.8 $\pm$ 5.3	54.6 $\pm$ 7	57.1 $\pm$ 6.5
-random-medium-expert	67.4 $\pm$ 11.9	15.5 $\pm$ 46.1	34 $\pm$ 5.1	95.3 $\pm$ 3.9	92.2 $\pm$ 2.9	95.6 $\pm$ 1.9
Maze ( <i>Actuators</i> = 5)						
-expert	100.1 $\pm$ 0	-1.7 $\pm$ 0	100.1 $\pm$ 0	100.1 $\pm$ 0	100.1 $\pm$ 0	100.1 $\pm$ 0
-medium-expert	99.6 $\pm$ 0.6	-1.7 $\pm$ 0	99.4 $\pm$ 0.9	100.1 $\pm$ 0	99 $\pm$ 2.3	100 $\pm$ 0.1
-medium	32.3 $\pm$ 4.3	-1.7 $\pm$ 0	32.7 $\pm$ 2.8	44.8 $\pm$ 8.4	44.8 $\pm$ 6.4	44.5 $\pm$ 3.6
-random-medium-expert	62.5 $\pm$ 12.1	-1.7 $\pm$ 0	21.4 $\pm$ 18	99.3 $\pm$ 1.4	99.4 $\pm$ 0.5	99.7 $\pm$ 0.4
Maze ( <i>Actuators</i> = 7)						
-expert	100.1 $\pm$ 0	-0.6 $\pm$ 0	100.1 $\pm$ 0	100.1 $\pm$ 0	100.1 $\pm$ 0	100.1 $\pm$ 0
-medium-expert	99.6 $\pm$ 0.6	-0.6 $\pm$ 0	98.4 $\pm$ 2.1	100.1 $\pm$ 0	100.1 $\pm$ 0	100.1 $\pm$ 0
-medium	68.1 $\pm$ 3	-0.6 $\pm$ 0	70.9 $\pm$ 4	76.8 $\pm$ 4.2	74 $\pm$ 4	73 $\pm$ 4.5
-random-medium-expert	83.3 $\pm$ 8.2	-0.6 $\pm$ 0	17 $\pm$ 12.5	99.2 $\pm$ 1.1	98.1 $\pm$ 1.5	99.2 $\pm$ 0.6
Maze ( <i>Actuators</i> = 10)						
-expert	100.1 $\pm$ 0	0 $\pm$ 0	100.1 $\pm$ 0	100.1 $\pm$ 0	100.1 $\pm$ 0	100.1 $\pm$ 0
-medium-expert	98.6 $\pm$ 1.2	0 $\pm$ 0	99.3 $\pm$ 0.8	100.1 $\pm$ 0	100 $\pm$ 0	100 $\pm$ 0
-medium	40.1 $\pm$ 5.8	0 $\pm$ 0	40 $\pm$ 6.4	53.4 $\pm$ 5.8	43.8 $\pm$ 5	42.2 $\pm$ 6.5
-random-medium-expert	84.5 $\pm$ 7	0 $\pm$ 0	50.9 $\pm$ 4.7	99.9 $\pm$ 0.5	97.8 $\pm$ 3	99.4 $\pm$ 0.7
Maze ( <i>Actuators</i> = 12)						
-expert	100.1 $\pm$ 0	$\pm$	$\pm$	$\pm$	$\pm$	$\pm$
-medium-expert	97.8 $\pm$ 1.6	0 $\pm$ 0	98.2 $\pm$ 1.5	100.1 $\pm$ 0	99 $\pm$ 0.7	99.8 $\pm$ 0.5
-medium	20.7 $\pm$ 3.2	0 $\pm$ 0	25.3 $\pm$ 6	38.1 $\pm$ 4.8	26.4 $\pm$ 4.2	30.1 $\pm$ 5.8
-random-medium-expert	79.7 $\pm$ 6.7	0 $\pm$ 0	65.3 $\pm$ 29.3	98.9 $\pm$ 1.1	95.7 $\pm$ 3.9	98.1 $\pm$ 1.9
Maze ( <i>Actuators</i> = 15)						
-expert	100.1 $\pm$ 0	0 $\pm$ 0	100.1 $\pm$ 0	100.1 $\pm$ 0	100.1 $\pm$ 0	100.1 $\pm$ 0
-medium-expert	100 $\pm$ 0	0 $\pm$ 0	99.7 $\pm$ 0.6	100.1 $\pm$ 0	100.1 $\pm$ 0	100.1 $\pm$ 0
-medium	53.2 $\pm$ 6.2	0 $\pm$ 0	50.8 $\pm$ 4.6	69.1 $\pm$ 3.4	91.3 $\pm$ 6.8	94.9 $\pm$ 2
-random-medium-expert	93.1 $\pm$ 3.2	0 $\pm$ 0	12.1 $\pm$ 17.9	99.3 $\pm$ 0.9	96.7 $\pm$ 4.9	85.1 $\pm$ 12.5
Sum						
-expert	600.6 $\pm$ 0	-7.9 $\pm$ 0	600.6 $\pm$ 0	600.7 $\pm$ 0	600.6 $\pm$ 0	600.6 $\pm$ 0
-medium-expert	594.3 $\pm$ 5	-7.9 $\pm$ 0	594.6 $\pm$ 6.5	600.6 $\pm$ 0	598.3 $\pm$ 3	599.9 $\pm$ 0.6
-medium	269.7 $\pm$ 25	-7.9 $\pm$ 0	274.4 $\pm$ 25.5	345 $\pm$ 31.9	334.9 $\pm$ 33.4	341.8 $\pm$ 28.9
-random-medium-expert	470.5 $\pm$ 49.1	13.2 $\pm$ 46.1	200.7 $\pm$ 87.5	591.9 $\pm$ 8.9	579.9 $\pm$ 16.7	577.1 $\pm$ 18
-all	1935.1 $\pm$ 79.1	-10.5 $\pm$ 46.1	1670.3 $\pm$ 119.5	2138.2 $\pm$ 40.8	2113.7 $\pm$ 53.1	2119.4 $\pm$ 47.5

Table 8: Individual performance comparison for DMC  $n = 3$ . Figures are mean normalised scores  $\pm$  one standard error, with 0 and 100 representing random and expert policies, respectively. At the bottom of the table we also provide totals, split by data quality and the entire set. We see that DecQN-BCQ is the least performant of the offline methods and that DecQN-CQL/IQL/OneStep perform similarly overall and on expert, medium-expert and and medium datasets and DecQN-CQL the best on random-medium-expert datasets.

Environment -dataset	BC	DecQN	DecQN-BCQ	DecQN-CQL	DecQN-IQL	DecQN-OneStep
FingerSpin						
-expert	99.5 $\pm$ 0.4	-0.2 $\pm$ 0.1	100.5 $\pm$ 0.8	107.1 $\pm$ 0.3	102.9 $\pm$ 0.2	102.5 $\pm$ 0.4
-medium-expert	77.9 $\pm$ 6.5	-0.5 $\pm$ 0	86.5 $\pm$ 9	106.8 $\pm$ 0.2	102.8 $\pm$ 0.3	102.7 $\pm$ 0.2
-medium	38.3 $\pm$ 1	-0.5 $\pm$ 0	40.7 $\pm$ 0.4	49.4 $\pm$ 0.3	44 $\pm$ 0.8	45.2 $\pm$ 0.6
-random-medium-expert	8.2 $\pm$ 2.3	56.2 $\pm$ 2.9	62.7 $\pm$ 9.5	100 $\pm$ 0.5	77.1 $\pm$ 5.1	78.5 $\pm$ 3.5
FishSwim						
-expert	82.9 $\pm$ 12.2	-1.8 $\pm$ 1.9	112 $\pm$ 2.7	120.8 $\pm$ 14.7	123.8 $\pm$ 11.7	105.3 $\pm$ 8.2
-medium-expert	40.6 $\pm$ 6.9	0.7 $\pm$ 6.5	97.5 $\pm$ 4.3	127.2 $\pm$ 11.1	91.1 $\pm$ 13.3	112.2 $\pm$ 8.5
-medium	42.8 $\pm$ 9.6	-4 $\pm$ 1.5	56.1 $\pm$ 9.8	71.5 $\pm$ 6.3	63.4 $\pm$ 7.2	76.7 $\pm$ 6
-random-medium-expert	23.6 $\pm$ 9.2	-0.4 $\pm$ 3.3	17.8 $\pm$ 4.8	52.1 $\pm$ 9	37.1 $\pm$ 9	59.9 $\pm$ 11
CheetahRun						
-expert	99.9 $\pm$ 1.7	2.7 $\pm$ 2.1	105.5 $\pm$ 0.4	105.6 $\pm$ 0.9	104.6 $\pm$ 1.1	106.3 $\pm$ 0.3
-medium-expert	61.6 $\pm$ 12.5	0.9 $\pm$ 0.8	104.2 $\pm$ 1.6	103.2 $\pm$ 0.7	102.5 $\pm$ 1.2	104.8 $\pm$ 0.7
-medium	40.4 $\pm$ 0.4	0 $\pm$ 0.7	47.1 $\pm$ 0.4	48.3 $\pm$ 0.3	47.7 $\pm$ 0.5	47.9 $\pm$ 0.3
-random-medium-expert	41.5 $\pm$ 0.6	0.5 $\pm$ 1.2	22.2 $\pm$ 4.4	79.6 $\pm$ 5.6	61.4 $\pm$ 3.3	53.2 $\pm$ 1.9
QuadrupedWalk						
-expert	97.7 $\pm$ 3.2	3.3 $\pm$ 4.6	114.9 $\pm$ 1.2	118.2 $\pm$ 1.4	122.3 $\pm$ 1.1	115.3 $\pm$ 2.4
-medium-expert	63.4 $\pm$ 11	8.5 $\pm$ 1.5	110.4 $\pm$ 1.9	115.4 $\pm$ 4	121.2 $\pm$ 1	109.4 $\pm$ 2.1
-medium	39.2 $\pm$ 8.5	11.3 $\pm$ 6.3	47.6 $\pm$ 8.4	48.6 $\pm$ 7.9	46.3 $\pm$ 5.7	46.8 $\pm$ 9.2
-random-medium-expert	28 $\pm$ 6.4	-5.2 $\pm$ 4.7	-12.3 $\pm$ 1.5	76.7 $\pm$ 5	65.8 $\pm$ 4.4	69.4 $\pm$ 5.7
HumanoidStand						
-expert	102.2 $\pm$ 1.3	-0.1 $\pm$ 0	103.6 $\pm$ 1.7	109 $\pm$ 1.9	116.6 $\pm$ 0.6	117.2 $\pm$ 0.8
-medium-expert	63.1 $\pm$ 3.9	0.1 $\pm$ 0	92 $\pm$ 2.7	104.7 $\pm$ 1.9	113.3 $\pm$ 0.4	116.7 $\pm$ 0.7
-medium	44.4 $\pm$ 0.5	0 $\pm$ 0	47.2 $\pm$ 1.5	51.4 $\pm$ 0.3	53.8 $\pm$ 0.4	53.8 $\pm$ 0.4
-random-medium-expert	34.4 $\pm$ 2.6	0 $\pm$ 0.1	17.7 $\pm$ 5.9	42.7 $\pm$ 0.9	46 $\pm$ 1	47.3 $\pm$ 0.7
DogTrot						
-expert	98 $\pm$ 0.7	0.1 $\pm$ 0.1	94.4 $\pm$ 0.5	99.5 $\pm$ 0.7	98.9 $\pm$ 2.2	101.2 $\pm$ 1.6
-medium-expert	62 $\pm$ 3.7	0 $\pm$ 0.3	74.9 $\pm$ 3.5	84.8 $\pm$ 3.7	89.3 $\pm$ 1.4	93.9 $\pm$ 1.4
-medium	43.8 $\pm$ 0.5	0.1 $\pm$ 0.1	49.1 $\pm$ 0.6	46.5 $\pm$ 0.5	52 $\pm$ 0.3	50.2 $\pm$ 0.2
-random-medium-expert	37.2 $\pm$ 3.6	0.1 $\pm$ 0.2	0.1 $\pm$ 0.1	43.4 $\pm$ 0.5	44.1 $\pm$ 1.2	44.9 $\pm$ 0.7
Sum						
-expert	580.2 $\pm$ 19.5	4 $\pm$ 8.8	630.9 $\pm$ 13.2	660.2 $\pm$ 19.9	669.1 $\pm$ 16.9	647.8 $\pm$ 13.7
-medium-expert	368.6 $\pm$ 44.5	9.7 $\pm$ 9.1	565.5 $\pm$ 36.9	642.1 $\pm$ 21.6	620.2 $\pm$ 17.6	639.7 $\pm$ 13.6
-medium	248.9 $\pm$ 20.5	6.9 $\pm$ 8.6	287.8 $\pm$ 44.8	315.7 $\pm$ 15.6	307.2 $\pm$ 14.9	320.6 $\pm$ 16.7
-random-medium-expert	172.9 $\pm$ 24.7	51.2 $\pm$ 12.4	108.2 $\pm$ 39.4	394.5 $\pm$ 21.5	331.5 $\pm$ 24	353.2 $\pm$ 23.5
-all	1370.6 $\pm$ 109.2	71.8 $\pm$ 38.9	1592.4 $\pm$ 134.3	2012.5 $\pm$ 78.6	1928 $\pm$ 73.4	1961.3 $\pm$ 67.5

Table 9: Individual performance comparison for dog-trot  $n = \{3, 10, 30, 50, 75, 100\}$ . Figures are mean normalised scores  $\pm$  one standard error, with 0 and 100 representing random and expert policies, respectively. At the bottom of the table we also provide totals, split by data quality and the entire set. We see across all datasets that DecQN-IQL/OneStep perform best, followed by DecQN-CQL and DecQN-BCQ.

Environment -dataset	BC	DecQN	DecQN-BCQ	DecQN-CQL	DecQN-IQL	DecQN-OneStep
DogTrot ( $n = 3$ )						
-expert	$98 \pm 0.7$	$0.1 \pm 0.1$	$94.4 \pm 0.5$	$99.5 \pm 0.7$	$98.9 \pm 2.2$	$101.2 \pm 1.6$
-medium-expert	$62 \pm 3.7$	$0 \pm 0.3$	$74.9 \pm 3.5$	$84.8 \pm 3.7$	$89.3 \pm 1.4$	$93.9 \pm 1.4$
-medium	$43.8 \pm 0.5$	$0.1 \pm 0.1$	$49.1 \pm 0.6$	$46.5 \pm 0.5$	$52 \pm 0.3$	$50.2 \pm 0.2$
-random-medium-expert	$37.2 \pm 3.6$	$0.1 \pm 0.2$	$0.1 \pm 0.1$	$43.4 \pm 0.5$	$44.1 \pm 1.2$	$44.9 \pm 0.7$
DogTrot ( $n = 10$ )						
-expert	$97.3 \pm 1.7$	$-0.3 \pm 0$	$96.8 \pm 1.8$	$99.2 \pm 1.3$	$106.5 \pm 1.1$	$113.4 \pm 0.9$
-medium-expert	$58.1 \pm 5.5$	$0.5 \pm 0.1$	$82.3 \pm 4.6$	$83.4 \pm 2.1$	$105.1 \pm 2.1$	$109.8 \pm 1.2$
-medium	$34.7 \pm 0.4$	$0.5 \pm 0.2$	$36.9 \pm 0.3$	$39.1 \pm 0.6$	$41 \pm 0.8$	$47.2 \pm 0.4$
-random-medium-expert	$33.6 \pm 4.4$	$0.1 \pm 0.1$	$6.1 \pm 2.8$	$33.4 \pm 1.5$	$52.1 \pm 2.3$	$45.4 \pm 4.1$
DogTrot ( $n = 30$ )						
-expert	$99.7 \pm 0.4$	$-0.3 \pm 0$	$100.3 \pm 2.3$	$99.8 \pm 0.8$	$102.8 \pm 0.4$	$106.4 \pm 0.8$
-medium-expert	$70.3 \pm 4.9$	$0.6 \pm 0.2$	$74.1 \pm 8$	$90.7 \pm 3.7$	$98.3 \pm 1.6$	$100.7 \pm 3$
-medium	$33.1 \pm 0.2$	$0.7 \pm 0.2$	$33.6 \pm 0.4$	$33.4 \pm 0.4$	$39.9 \pm 0.1$	$40.5 \pm 0.5$
-random-medium-expert	$22.4 \pm 0.9$	$0.2 \pm 0.1$	$5.2 \pm 2.2$	$22.9 \pm 2.5$	$31.6 \pm 1.3$	$30.4 \pm 1.8$
DogTrot ( $n = 50$ )						
-expert	$98.2 \pm 0.5$	$0.5 \pm 0.1$	$98.6 \pm 0.6$	$97.7 \pm 0.6$	$99.2 \pm 0.9$	$100.4 \pm 1.4$
-medium-expert	$67.7 \pm 5$	$-0.3 \pm 0$	$68.7 \pm 4.2$	$88.9 \pm 2.1$	$95 \pm 2.1$	$97.1 \pm 1.4$
-medium	$34.9 \pm 0.4$	$-0.3 \pm 0$	$36.5 \pm 0.7$	$36.7 \pm 0.2$	$45.1 \pm 0.6$	$43.1 \pm 0.7$
-random-medium-expert	$21 \pm 2$	$0.7 \pm 0.1$	$9.3 \pm 1.5$	$27.6 \pm 1.3$	$27.9 \pm 2.4$	$36.4 \pm 1.8$
DogTrot ( $n = 75$ )						
-expert	$98.4 \pm 0.8$	$0.6 \pm 0.2$	$99.6 \pm 1.1$	$98 \pm 1$	$101.4 \pm 1.4$	$102.9 \pm 0.9$
-medium-expert	$75.3 \pm 2.6$	$0 \pm 0$	$74.9 \pm 16.7$	$85 \pm 2$	$89.9 \pm 0.7$	$97.9 \pm 2.5$
-medium	$41.4 \pm 0.2$	$0.5 \pm 0.1$	$41.3 \pm 0.8$	$42.5 \pm 0.1$	$48.5 \pm 0.4$	$47.8 \pm 0.5$
-random-medium-expert	$27.7 \pm 0.7$	$0.7 \pm 0.2$	$0 \pm 0.1$	$33.4 \pm 2$	$45.3 \pm 2.9$	$43.6 \pm 2.6$
DogTrot ( $n = 100$ )						
-expert	$96.4 \pm 1.1$	$0.1 \pm 0$	$95 \pm 1.8$	$100.4 \pm 0.7$	$103.2 \pm 0.4$	$105.5 \pm 0.8$
-medium-expert	$70.9 \pm 6.5$	$0.4 \pm 0.1$	$74.6 \pm 5.8$	$74.4 \pm 5.1$	$93.9 \pm 1$	$101.5 \pm 2.5$
-medium	$34.7 \pm 0.8$	$0.1 \pm 0.1$	$34.5 \pm 0.5$	$34.8 \pm 0.7$	$40.9 \pm 0.2$	$41.5 \pm 0.5$
-random-medium-expert	$18.9 \pm 2.7$	$0.3 \pm 0.1$	$0 \pm 0$	$12.7 \pm 2.4$	$24 \pm 2.8$	$26.2 \pm 1$
Sum						
-expert	$588 \pm 5.2$	$0.7 \pm 0.4$	$584.7 \pm 14$	$594.6 \pm 5.1$	$612 \pm 6.4$	$629.8 \pm 6.4$
-medium-expert	$404.3 \pm 28.2$	$1.2 \pm 0.7$	$449.5 \pm 64.9$	$507.2 \pm 18.7$	$571.5 \pm 8.9$	$600.9 \pm 12$
-medium	$222.6 \pm 2.5$	$1.6 \pm 0.7$	$231.9 \pm 5.8$	$233 \pm 2.5$	$267.4 \pm 2.4$	$270.3 \pm 2.8$
-random-medium-expert	$160.8 \pm 14.3$	$2.1 \pm 0.8$	$30.2 \pm 17.6$	$173.4 \pm 10.2$	$225 \pm 12.9$	$226.9 \pm 12$
-all	$1375.7 \pm 50.2$	$5.6 \pm 2.6$	$1296.3 \pm 102.3$	$1508.2 \pm 36.5$	$1675.9 \pm 30.6$	$1727.9 \pm 33.2$

Table 10: Performance across 5 seeds.

Environment -dataset	DecQN-BCQ	DecQN-CQL	DecQN-IQL	DecQN-OneStep
HumanoidStand				
-expert	$103.6 \pm 1.7$	$109 \pm 1.9$	$116.6 \pm 0.6$	$117.2 \pm 0.8$
-medium-expert	$92 \pm 2.7$	$104.7 \pm 1.9$	$113.3 \pm 0.4$	$116.7 \pm 0.7$
-medium	$47.2 \pm 1.5$	$51.4 \pm 0.3$	$53.8 \pm 0.4$	$53.8 \pm 0.4$
-random-medium-expert	$17.7 \pm 5.9$	$42.7 \pm 0.9$	$46 \pm 1$	$47.3 \pm 0.7$



Table 11: Performance across 10 seeds

Environment -dataset	DecQN-BCQ	DecQN-CQL	DecQN-IQL	DecQN-OneStep
HumanoidStand				
-expert	$103.6 \pm 1.3$	$109.6 \pm 1.5$	$117 \pm 0.6$	$118.1 \pm 0.7$
-medium-expert	$90.3 \pm 3.3$	$103.4 \pm 1.8$	$112.9 \pm 1.1$	$116.2 \pm 0.6$
-medium	$46.7 \pm 0.6$	$51.9 \pm 0.4$	$53.7 \pm 0.4$	$53.8 \pm 0.4$
-random-medium-expert	$19.6 \pm 2.5$	$41.4 \pm 1.1$	$45.5 \pm 1.3$	$46.9 \pm 0.5$

Table 12: Individual performance comparison for DecQN-CQL using mean and sum decompositions and BDQ-CQL for  $n = 3$ . Figures are mean normalised scores  $\pm$  one standard error, with 0 and 100 representing random and expert policies, respectively.

Environment -dataset	DecQN-CQL (Mean)	DecQN-CQL (Sum)	BDQ-CQL
HumanoidStand			
-expert	$109 \pm 1.9$	$95.1 \pm 1.8$	$105.5 \pm 0.7$
-medium-expert	$104.7 \pm 1.9$	$86.1 \pm 2.8$	$94.0 \pm 5.8$
-medium	$51.4 \pm 0.3$	$44.6 \pm 0.8$	$48.2 \pm 0.5$
-random-medium-expert	$42.7 \pm 0.9$	$36.3 \pm 1.6$	$43.2 \pm 0.5$
DogTrot			
-expert	$99.5 \pm 0.7$	$93.1 \pm 1.3$	$94.9 \pm 1.2$
-medium-expert	$84.8 \pm 3.7$	$81.4 \pm 2.8$	$75.7 \pm 3.6$
-medium	$46.5 \pm 0.5$	$44.3 \pm 0.3$	$48.5 \pm 0.5$
-random-medium-expert	$43.4 \pm 0.5$	$39.9 \pm 0.9$	$37.9 \pm 2.3$