

Adaptive Trade-Offs in Off-Policy Learning

Mark Rowland*

DeepMind

Will Dabney*

DeepMind

Rémi Munos

DeepMind

Abstract

A great variety of off-policy learning algorithms exist in the literature, and new breakthroughs in this area continue to be made, improving theoretical understanding and yielding state-of-the-art reinforcement learning algorithms. In this paper, we take a unifying view of this space of algorithms, and consider their trade-offs of three fundamental quantities: *update variance*, *fixed-point bias*, and *contraction rate*. This leads to new perspectives of existing methods, and also naturally yields novel algorithms for off-policy evaluation and control. We develop one such algorithm, C-trace, demonstrating that it is able to more efficiently make these trade-offs than existing methods in use, and that it can be scaled to yield state-of-the-art performance in large-scale environments.

1 Introduction

Off-policy learning is crucial to modern reinforcement learning, allowing agents to learn from memorised data, demonstrations, and exploratory behaviour [Szepesvári, 2010, Sutton and Barto, 2018]. As such, it is a long-studied problem, with a variety of well-understood associated algorithms; see [Precup et al., 2000, Kakade and Langford, 2002, Dudík et al., 2014, Thomas and Brunskill, 2016, Munos et al., 2016, Mahmood et al., 2017, Farajtabar et al., 2018] for a representative selection of publications.

However, this paper is motivated by the observation that in spite of this theoretical progress, several state-of-the-art value-based reinforcement learning agents (notably Rainbow [Hessel et al., 2018] and R2D2 [Kapturowski et al., 2019]) eschew these off-policy algorithms, attaining better performance by using *uncorrected* re-

turns, which do not account for the fact that data is generated off-policy. Further research has corroborated this observation [Hernandez-Garcia and Sutton, 2019]. This raises two central research questions: (i) How can we understand the strong performance of uncorrected returns? (ii) Can we distil these advantages, and combine them with existing off-policy algorithms to improve their performance?

One of the principal contributions of this paper is to show that the performance of all off-policy evaluation algorithms can be decomposed into just three fundamental quantities: *contraction rate*, *fixed-point bias*, and *variance*; see Figure 1 for a preliminary illustration. Intuitively, *fixed-point bias* describes the error of an algorithm in the limit of infinite data, *contraction rate* describes the speed at which an algorithm approaches its infinite-data limit, and *variance* describes to what extent randomly observed data can perturb the algorithm.

This decomposition yields an interpretation of the empirical success of uncorrected returns, and an answer to question (i) above; namely, that they are efficiently making a trade-off between fixed-point bias and the other fundamental quantities. Further, this suggests an answer to question (ii) — that we may be able to improve existing off-policy algorithms by incorporating a means of making such a trade-off. This leads us to the development of *C-trace*, a new off-policy algorithm that achieves strong empirical performance in several large-scale environments.

We develop the trade-off framework mentioned above in Section 2, proving the existence of the three fundamental quantities described above, and showing that all off-policy algorithms necessarily make an implicit trade-off between these quantities. We then use this framework to develop a new off-policy learning algorithm, C-trace, in Section 3, and study its contraction and convergence properties. We then demonstrate its empirical effectiveness in tabular domains and when applied to two deep reinforcement learning agents, DQN [Mnih et al., 2015] and R2D2 [Kapturowski et al., 2019], in Section 4.

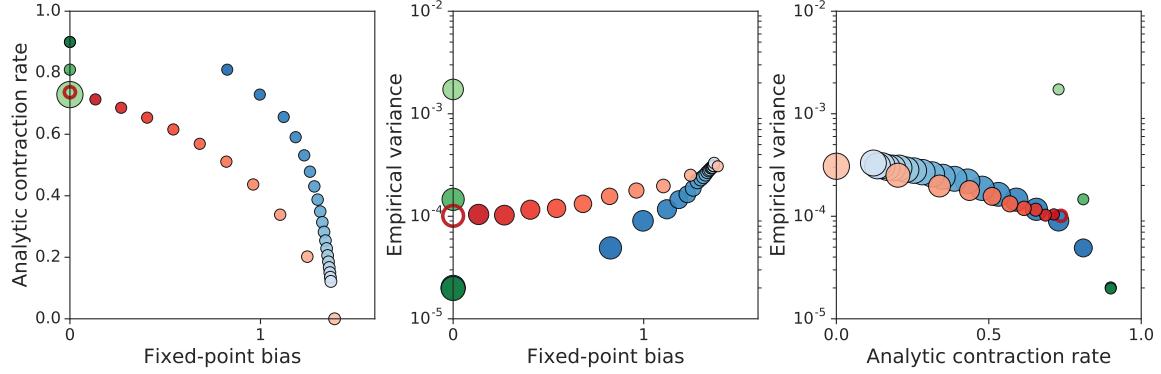


Figure 1: Trade-offs made by n -step uncorrected returns (dark blue [$n = 1$] through to light blue [$n = 20$]), n -step importance corrected returns (dark green [$n = 1$] through to light green [$n = 3$]), Retrace (red open circle). Also pictured is the new method α -Retrace (dark red [$\alpha = 1$] through to light red [$\alpha = 0$]), introduced in Section 3. All quantities are calculated for a fixed evaluation problem in a small, randomly generated MDP; see Appendix Section C.1 for further environment details. In each plot, the magnitude of the points illustrates the relative scale of the third trade-off variable.

1.1 Notation and preliminary definitions

Throughout, we consider a Markov decision process (MDP) with finite state space \mathcal{X} , finite action space \mathcal{A} , discount factor $\gamma \in [0, 1)$, transition kernel $P : \mathcal{X} \times \mathcal{A} \rightarrow \mathcal{P}(\mathcal{X})$, reward distributions $\mathcal{R} : \mathcal{X} \times \mathcal{A} \rightarrow \mathcal{P}(\mathbb{R})$, and some initial state distribution $\nu_0 \in \mathcal{P}(\mathcal{X})$. Given a Markov policy $\pi : \mathcal{X} \rightarrow \mathcal{P}(\mathcal{A})$, we write $(X_t, A_t, R_t)_{t \geq 0}$ for the process describing the sequence of states visited, actions taken, and rewards received by an agent acting in the MDP according to π , so that $R_t | X_t, A_t \sim \mathcal{R}(X_t, A_t)$ for all $t \geq 0$. Additionally, we write $r(x, a)$ for the expected immediate reward received after taking action a in state x . Given a policy π , the task of *evaluation* is to learn the function $Q^\pi(x, a) = \mathbb{E}_\pi [\sum_{t=0}^{\infty} \gamma^t R_t | X_0 = x, A_0 = a]$, where \mathbb{E}_π denotes expectation with respect to the distribution over trajectories induced by π . The task of *control* is to identify the Markov policy π^* maximising the quantity $\mathbb{E}[Q^\pi(X_0, A_0)]$, where $A_0 \sim \pi(\cdot | X_0)$, and $X_0 \sim \nu_0$. We also define the one-step evaluation operator $T^\mu : \mathbb{R}^{\mathcal{X} \times \mathcal{A}} \rightarrow \mathbb{R}^{\mathcal{X} \times \mathcal{A}}$ associated with a Markov policy $\mu : \mathcal{X} \rightarrow \mathcal{P}(\mathcal{A})$ by

$$(T^\mu Q)(x, a) = r(x, a) + \sum_{x' \in \mathcal{X}, a' \in \mathcal{A}} P(x'|x, a)\mu(a'|x')Q(x', a'), \quad (1)$$

for all $Q \in \mathbb{R}^{\mathcal{X} \times \mathcal{A}}$, and $(x, a) \in \mathcal{X} \times \mathcal{A}$.

We now briefly give formal definitions of the key concepts we seek to analyse in this paper.

Definition 1.1. An **evaluation update rule** for evaluating a policy π under a behaviour policy μ is a function $\hat{T} : \mathbb{R}^{\mathcal{X} \times \mathcal{A}} \times (\mathcal{X} \times \mathcal{A} \times \mathbb{R})^* \rightarrow \mathbb{R}$ that takes

as input a value function estimate Q and a trajectory $(x_t, a_t, r_t)_{t \geq 0}$ given by following μ , and outputs an update for $Q(x_0, a_0)$. There is an associated **evaluation operator** $T : \mathbb{R}^{\mathcal{X} \times \mathcal{A}} \rightarrow \mathbb{R}^{\mathcal{X} \times \mathcal{A}}$, given by

$$(TQ)(x, a) = \mathbb{E}_\mu \left[\hat{T}(Q, (X_t, A_t, R_t)_{t=0}^{\infty}) \middle| X_0 = x, A_0 = a \right],$$

for all $Q \in \mathbb{R}^{\mathcal{X} \times \mathcal{A}}$ and $(x, a) \in \mathcal{X} \times \mathcal{A}$.

Definition 1.2. The **contraction rate** of an operator $T : \mathbb{R}^{\mathcal{X} \times \mathcal{A}} \rightarrow \mathbb{R}^{\mathcal{X} \times \mathcal{A}}$ is given by

$$\Gamma = \sup_{\substack{Q, Q' \in \mathbb{R}^{\mathcal{X} \times \mathcal{A}} \\ Q \neq Q'}} \|TQ - TQ'\|_\infty / \|Q - Q'\|_\infty,$$

An operator is said to be *contractive* if $\Gamma < 1$. We can also consider state-action specific contraction rates via the quantities $\sup_{Q \neq Q'} |(TQ)(x, a) - (TQ')(x, a)| / \|Q - Q'\|_\infty$.

Definition 1.3. For a contractive operator T targeting a policy π , the **fixed-point bias** of T is given by $\|Q^\pi - \hat{Q}^\pi\|_2$, where \hat{Q}^π is the unique fixed point of T (guaranteed to exist by the contractivity of T).

Definition 1.4. The **variance** of an update rule \hat{T} stochastically approximating an operator T at approximate value function Q and an initial state-action distribution $\nu \in \mathcal{P}(\mathcal{X} \times \mathcal{A})$ is $\mathbb{E}_{(X_0, A_0) \sim \nu} \left[\mathbb{E}_\mu \left[\|\hat{T}(Q, (X_t, A_t, R_t)_{t=0}^{\infty}) - TQ\|_2^2 \middle| X_0, A_0 \right] \right]$.

2 Contraction, bias, and variance in off-policy evaluation

We begin with two motivating examples from recent research in off-policy evaluation methods, illustrating

examples of the types of trade-offs we seek to describe in this paper.

***n*-step uncorrected returns.** Recently proposed agents such as Rainbow [Hessel et al., 2018] and R2D2 [Kapturowski et al., 2019] have made use of the *uncorrected n-step return* in constructing off-policy learning algorithms. Consistent with these results, Hernandez-Garcia and Sutton [2019] observed that these uncorrected updates frequently outperformed off-policy corrections. Given an estimate \hat{Q} of the action-value function Q^π , the *n*-step uncorrected target for $\hat{Q}(x_0, a_0)$, given a trajectory $(x_0, a_0, r_0, x_1, a_1, r_1, \dots, x_n)$ of experience generated according to behaviour policy μ , is given by

$$\sum_{s=0}^{n-1} \gamma^s r_s + \gamma^n \mathbb{E}_{A \sim \pi(\cdot|x_n)} [\hat{Q}(x_n, A)]. \quad (2)$$

The adjective *uncorrected* contrasts this update target against the *n-step importance-weighted return* target, which takes the following form:

$$\sum_{s=0}^{n-1} \rho_{1:s} \gamma^s r_s + \rho_{1:n-1} \gamma^n \mathbb{E}_{A \sim \pi(\cdot|x_n)} [\hat{Q}(x_n, A)], \quad (3)$$

where we write $\rho_t = \pi(a_t|x_t)/\mu(a_t|x_t)$, and $\rho_{s:t} = \prod_{u=s}^t \rho_u$ for each $1 \leq s \leq t$. Empirically, the former has been observed to work very well in these recent works, whilst the latter is often too unstable to be used; this fact is often attributed to the *high variance* of the importance-weighted update, with the uncorrected update having relatively low variance by comparison. We also observe that the uncorrected update is a stochastic approximation to the operator $(T^\mu)^{n-1} T^\pi$, whilst the importance-weighted update is a stochastic approximation to $(T^\pi)^n$. From this, it follows that under usual stochastic approximation conditions, a sequence of importance-weighted updates will converge to the true action-function Q^π associated with π , whilst the uncorrected updates will converge to the value function of the time-inhomogeneous policy that follows π for a single step, followed by $n - 1$ steps of μ , and then repeats; see Proposition B.1 in Appendix Section B for further explanation.

The above discussion shows that we may view the use of uncorrected returns as trading off update *variance* for *accuracy of the operator fixed point*; an example of the classical bias-variance trade-off in statistics and machine learning, albeit in the context of fixed-point iteration.

Retrace. Munos et al. [2016] proposed an off-policy evaluation update target, Retrace, given in its forward-view version by

$$\hat{Q}(x_0, a_0) + \sum_{s \geq 0} \bar{\rho}_{1:s} \gamma^s \Delta_s, \quad (4)$$

where we write $\bar{\rho}_t = \min(1, \rho_t)$, and $\bar{\rho}_{s:t} = \prod_{u=s}^t \bar{\rho}_u$ for each $1 \leq s \leq t$, and define the temporal difference (TD) error Δ_s by

$$\Delta_s \stackrel{\text{def}}{=} r_s + \gamma \mathbb{E}_{A \sim \pi(\cdot|x_{s+1})} [\hat{Q}(x_{s+1}, A)] - \hat{Q}(x_s, a_s).$$

By clipping the importance weights associated with each TD error, the variance associated with the update rule is reduced relative to importance-weighted returns, whilst no bias is introduced; the fixed point of the associated Retrace operator remains the true action-value function Q^π . However, the clipping of the importance weights effectively *cuts the traces* in the update, resulting in the update placing less weight on later TD errors, and thus worsening the contraction rate of the corresponding operator. Thus, Retrace can be interpreted as trading off a *reduction in update variance* for a *larger contraction rate*, relative to importance-weighted *n*-step returns.

We discuss more examples of off-policy learning algorithms in Section 5. We also note that λ -variants of the algorithms described above also exist; for clarity and conciseness, we limit our exposition to the case $\lambda = 1$ in the main paper, noting that the results straightforwardly extend to $\lambda \in (0, 1)$.

We now briefly return to Figure 1, which quantitatively illustrates the trade-offs discussed above. We highlight several interesting observations. Whilst all importance-weighted updates have no fixed-point bias, their variance grows exceptionally quickly with n . Retrace manages to achieve a similar contraction rate to the 3-step importance-weighted update, but without incurring high variance. Our new algorithm, α -Retrace, appears to be Pareto efficient relative to the *n*-step uncorrected methods in the left-most plot; for any contraction rate that an *n*-step uncorrected method achieves, there is a value of α such that α -Retrace achieves this contraction rate whilst incurring less fixed-point bias; this is corroborated by further empirical results in Appendix Section B.

2.1 Downstream tasks and bounds

Whilst the trade-offs at the level of individual updates described above are straightforward to describe, in reinforcement learning we are ultimately interested in one of two problems, either *evaluation* or *control*, defined formally below.

The evaluation problem. Given a target policy π , a budget of experience generated from a behaviour policy μ , and a computational budget, compute an accurate approximation \hat{Q} to Q^π , in the sense of incurring low error $\|\hat{Q} - Q^\pi\|$, for some norm $\|\cdot\|$.

The control problem. Given a budget of experience and computation, find a policy π such that expected return under π is maximised.

It is intuitively clear that for each of these problems, an evaluation scheme with low contraction rate, low update variance, and low fixed-point bias is advantageous, but no update is known to possess all three of these attributes simultaneously. What is less clear is how these three properties should be traded off against one another in designing an efficient off-policy learning algorithm. For example, how much fixed-point accuracy should one be willing trade off in exchange for a halved update variance? Questions such as these in general have complicated dependence on the precise structure of the update rule, the policies in question, and the environment too, and so it seems unlikely that much progress can be made here in great generality. However, it is possible to make some progress.

Proposition 2.1. Consider the task of evaluation of a policy π under behaviour μ , and consider an update rule \hat{T} which stochastically approximates the application of an operator \tilde{T} , with contraction rate Γ and fixed point \tilde{Q} , to an initial estimate Q . Then we have the following decomposition:

$$\mathbb{E} \left[\|\hat{T}Q - Q^\pi\|_\infty \right] \leq \underbrace{\mathbb{E} \left[\|\hat{T}Q - \tilde{T}Q\|_2^2 \right]^{1/2}}_{(\text{Root variance})} + \underbrace{\Gamma \|Q - \tilde{Q}\|_\infty}_{(\text{Contraction})} + \underbrace{\|\tilde{Q} - Q^\pi\|_2}_{(\text{Fixed-point bias})}.$$

This result gives some sense of how these trade-offs feed into evaluation quality; related decompositions are also possible, which we describe in Appendix Section B. We next show that there really is a trade-off to be made, in the sense that it is not possible for an update based on limited data to simultaneously have low variance, contraction rate, and fixed-point bias across a range of MDPs.

Theorem 2.2. Consider an update rule \hat{T} with corresponding operator T , and consider the collection $\mathcal{M} = \mathcal{M}(\mathcal{X}, \mathcal{A}, P, \gamma, R_{\max})$ of MDPs with common state space, action space, transition kernel, and discount factor (but varying rewards, with maximum immediate reward bounded in absolute value by R_{\max}). Fix a target policy π , and a random variable Z , the set of transitions used by the operator \hat{T} ; these could be transitions encountered in a trajectory following the behaviour μ , or i.i.d. samples from the discounted state-action visitation distribution under μ . We denote the mismatch between π and Z at level $\delta \in (0, 1)$ by

$$D(Z, \pi, \delta) \stackrel{\text{def}}{=} \max \{d_{(x,a),\pi}(\Omega) \mid \Omega \subseteq \mathcal{X} \times \mathcal{A} \text{ s.t. } \mathbb{P}(Z \cap \Omega \neq \emptyset) \leq \delta, (x, a) \in \mathcal{X} \times \mathcal{A}\},$$

where $d_{(x,a),\pi}$ is the discounted state-action visitation distribution for trajectories initialised with (x, a) , following π thereafter. Denoting the variance, contraction rate, and fixed-point bias of \hat{T} for a particular MDP $M \in \mathcal{M}$ by $\mathbb{V}(M)$, $\Gamma(M)$ and $B(M)$ respectively, we have

$$\sup_{M \in \mathcal{M}} \left[\sqrt{\mathbb{V}(M)} + \frac{2R_{\max}}{1-\gamma} \Gamma(M) + B(M) \right] \geq \sup_{\delta \in (0,1)} (1-\delta) D(Z, \pi, \delta) R_{\max} / (1-\gamma).$$

In addition to the above results, which we believe to be novel, there is extensive literature exploring particular aspects of these trade-offs, which we discuss further in Section 5. Having made this space of trade-offs between contraction, bias, and variance explicit, a natural questions is how other update rules might be modified to exploit different parts of the space. In particular, incurring some amount of fixed-point bias for reduced variance made by n -step uncorrected returns in Rainbow and R2D2 is particularly effective in practice — is there a way to introduce a similar trade-off in an algorithm with adaptive trace lengths, such as Retrace? We explore this question in the next section.

3 New off-policy updates: α -Retrace and C-trace

The Retrace update in Equation (4) has been observed, in certain scenarios, to cut traces prematurely [Mahmood et al., 2017]; that is, using n -step uncorrected returns for suitable n leads to a superior contraction rate relative to Retrace, outweighing the corresponding incurred bias. A natural question is how Retrace can be modified to overcome this phenomenon. In the language of Section 2, is there a way that Retrace can be adapted so as to trade off contraction rate for fixed-point bias? The reduced contraction rate comes from cases where the truncated importance weights $\min(1, \pi(a_t|x_t)/\mu(a_t|x_t))$ appearing in (4) are small, so a natural way to improve the contraction rate is to move the target policy closer towards the behaviour.

Algorithm 1 α -Retrace for policy iteration

```

Initialise target policy  $\tilde{\pi}$  and behaviour  $\mu$ .
for each policy improvement round: do
    Select  $\alpha \in [0, 1]$ , and set new target policy  $\pi = \alpha\tilde{\pi} + (1-\alpha)\mu$ .
    Learn  $\hat{Q}^\pi : \mathcal{X} \times \mathcal{A} \rightarrow \mathbb{R}$  via Retrace under behaviour policy  $\mu$ .
    Set  $\tilde{\pi} = \text{Greedy}(\hat{Q}^\pi)$ .
    Set new behaviour policy  $\mu$ .
end for

```

To this end, we propose α -Retrace, a family of algorithms that applies Retrace to a target policy given by a mixture of the original target and the behaviour, thus achieving the aforementioned trade-off. In Algorithm 1, we describe how α -Retrace can be used within a (modified) policy iteration scheme for control. Note that 1-Retrace is simply the standard Retrace algorithm. We refer back to Figure 1, the left-most plot of which shows that this mixture coefficient precisely yields a trade-off between fixed-point bias and contraction rate that we sought at the end of Section 2.

The means by which α should be set is left open at this stage; adjusting it allows a trade-off of contraction rate and fixed-point bias. In Section 3.2, we describe a stochastic approximation procedure for updating α online to obtain a desired contraction rate.

Specificity to Retrace. Whilst the mixture target of α -Retrace is a natural choice, we highlight that this choice is in fact specific to the structure of Retrace. In Appendix Section D.1, we visualise trade-offs made by analogous adjustment to the TreeBackup update [Precup et al., 2000], showing that mixing the behaviour policy into the target simply leads to an accumulation of fixed-point bias, with limited benefits in terms of contraction rate or variance.

3.1 Analysis

We now provide several results describing the contraction rate of α -Retrace in more detail, and how the fixed-point bias introduced by $\alpha < 1$ may be useful in the case of control tasks. We begin with a preliminary definition.

Definition 3.1. For a state-action pair $(x, a) \in \mathcal{X} \times \mathcal{A}$, Two policies π_1, π_2 are said to be (x, a) -**distinguishable** under a third policy μ if there exists $x' \in \mathcal{X}$ in the support of the discounted state visitation distribution under μ starting from state-action pair (x, a) , such that $\pi_1(\cdot|x') \neq \pi_2(\cdot|x')$, and are said to be (x, a) -**indistinguishable** under μ otherwise.

Proposition 3.2. The operator associated with the α -Retrace evaluation update for evaluating π given behaviour μ has a state-action-dependent contraction rate of

$$C(\alpha|x, a) \stackrel{\text{def}}{=} 1 - (1 - \gamma) \times \mathbb{E}_{\mu} \left[\sum_{t=0}^{\infty} \gamma^t \prod_{s=1}^t ((1 - \alpha) + \alpha \bar{\rho}_s) \mid (X_0, A_0) = (x, a) \right], \quad (5)$$

for each $(x, a) \in \mathcal{X} \times \mathcal{A}$. Viewed as a function of $\alpha \in [0, 1]$, this contraction rate is continuous, monotonically increasing, with minimal value 0, and maximal value no greater than γ . Further, the contraction rate is

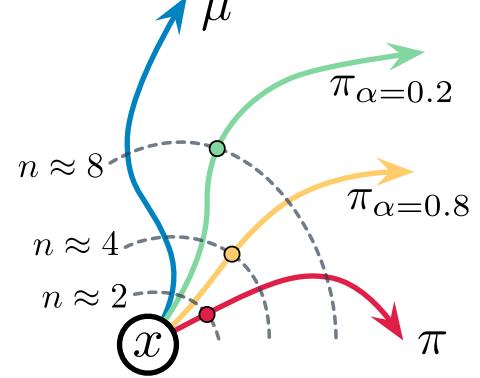


Figure 2: Interpolating between target policy π and behaviour policy μ with $\alpha \in \{0.0, 0.2, 0.8, 1.0\}$ produces different expected trajectories shown by each coloured line. As the mixture policy more closely resembles the behaviour policy, α -Retrace allows more off-policy data to be used (dashed line, numbers indicate expected trace-length), cuts traces (coloured points) later, yielding lower contraction rates equivalent to n -step methods with larger n . C-trace adapts α online to achieve a stable trace length throughout training.

strictly monotonic iff π and μ are (x, a) -distinguishable under μ .

The exact contraction rate of α -Retrace is thus $\sup_{(x, a) \in \mathcal{X} \times \mathcal{A}} C(\alpha|x, a)$, which inherits the continuity and monotonicity properties of the state-action-dependent rates. Our next result motivates the use of α -Retrace within control algorithms.

Proposition 3.3. Consider a target policy π , let μ be the behavioural policy, and assume that there is a unique greedy action $a^*(x) \in \mathcal{A}$ with respect to Q^π at each state x for each $x \in \mathcal{X}$. Then there exists a value of $\alpha \in (0, 1)$ such that the greedy policy with respect to the fixed point of α -Retrace coincides with the greedy policy with respect to Q^π , and the contraction rate for this α -Retrace is no greater than that for 1-Retrace. Further, if π and μ are (x, a) -distinguishable under μ for all $(x, a) \in \mathcal{X} \times \mathcal{A}$, then the contraction rate of α -Retrace is strictly lower than that of 1-Retrace.

3.2 C-trace: adapting α online

An empirical shortcoming of Retrace noted earlier is its tendency to pessimistically cut traces. Adapting the mixture parameter α within Retrace(α) yields a natural way to ensure that a desired trace length (or contraction rate) is attained. In this section, we propose C-trace, which uses α -Retrace updates whilst dynamically adjusting α to attain a target contraction rate Γ ; a schematic illustration is given in Figure 2.

The contraction rate $\sup_{(x, a) \in \mathcal{X} \times \mathcal{A}} C(\alpha|x, a)$ is difficult

to estimate online, so we work instead with the averaged contraction rate $C_\nu(\alpha) = \mathbb{E}_{(X,A) \sim \nu}[C(\alpha|X, A)]$, where ν is the training distribution over state-action pairs; where clear, we will drop ν from notation. It follows straightforwardly from Proposition 3.2 that $C_\nu(\alpha)$ is monotonic in α . This suggests that a standard Robbins-Monro stochastic approximation update rule for α may be applied to guide $C_\nu(\alpha)$ towards Γ — we describe such a scheme below. To avoid optimisation issues with the constraint $\alpha \in [0, 1]$, we parameterise α as $\sigma(\phi)$, where σ is the standard sigmoid function, and $\phi \in \mathbb{R}$ is an unconstrained real variable. For brevity, we will simply write $\alpha(\phi)$. Since σ is monotonic and continuous, the contraction rate is still monotonic and continuous in ϕ .

Recall from (5) that the contraction rate $C(\alpha|x, a)$ of the α -Retrace operator with target π and behaviour μ can be expressed as an expectation over trajectories following μ , and thus can be unbiasedly approximated using such trajectories; given an i.i.d. sequence of trajectories $(x_t^{(k)}, a_t^{(k)}, r_t^{(k)})_{t \geq 0}$, we write $\hat{C}^{(k)}(\alpha(\phi))$ for the corresponding estimates of $C(\alpha(\phi))$. If the target contraction rate is Γ , we can adjust an initial parameter $\phi_0 \in \mathbb{R}$ using these estimates according to the Robbins-Monro rule

$$\phi_{k+1} = \phi_k - \varepsilon_k \left(\hat{C}^{(k)}(\alpha(\phi_k)) - \Gamma \right) \quad \forall k \geq 0, \quad (6)$$

for some sequence of stepsizes $(\varepsilon_k)_{k=0}^\infty$. The following result gives a theoretical guarantee for the correctness of this procedure.

Proposition 3.4. Let $(x_t^{(k)}, a_t^{(k)}, r_t^{(k)})_{t \geq 0}^\infty$ be an i.i.d. sequence of trajectories following μ , with initial state-action distribution given by ν . Let Γ be a target contraction rate such that $C_\nu(1) \geq \Gamma$. Let the stepsizes $(\varepsilon_k)_{k=0}^\infty$ satisfy the usual Robbins-Monro conditions $\sum_{k=0}^\infty \varepsilon_k = \infty$, $\sum_{k=0}^\infty \varepsilon_k^2 < \infty$. Then for any initial value ϕ_0 following the updates in (6), we have $\phi_k \rightarrow \phi^*$ in probability, where $\phi^* \in \mathbb{R}$ is the unique value such that $C_\nu(\alpha(\phi^*)) = \Gamma$.

C-trace thus consists of interleaving α -Retrace evaluation updates with α parameter updates as in (6).

Convergence analysis. It is possible to further develop the theory in Proposition 3.4 to prove convergence of C-trace as a whole, using techniques going back to those of Bertsekas and Tsitsiklis [1996] for convergence of TD(λ), and more recently used by Munos et al. [2016] to prove convergence of a control version of Retrace, as the following result shows.

Theorem 3.5. Assume the same conditions as Proposition 3.4, and additionally that: (i) trajectory lengths have finite second moment; (ii) immediate rewards are bounded. Let $(\phi_k)_{k=0}^\infty$ be defined as in Equation (6) and $(Q_k)_{k=0}^\infty$ be a sequence of Q-functions, with

Q_{k+1} obtained from applying Retrace updates targeting $\alpha(\phi_k)\pi + (1 - \alpha(\phi_k))\mu$ to Q_k with trajectory $k+1$, using stepsize ε_k . Then we have $\alpha(\phi_k) \rightarrow \alpha(\phi^*) =: \alpha^*$ and $Q_k \rightarrow Q^{\alpha^*}\pi^{+(1-\alpha^*)\mu}$ almost surely.

Truncated trajectory corrections. The method described above for adaptation of α is impractical in scenarios where episodes are particularly long, when the MDP is non-episodic, and when only partial segments of trajectories can be processed at once. Since such cases often arise in practice, this motivates modifications to the update of (6). Here, we describe one such modification which will be crucial to the deployment of C-trace in large-scale agents in Section 4. Given a *truncated trajectory* $(x_t, a_t, r_t)_{t=0}^N$, Retrace necessarily must cut traces after at most N time steps, and so can achieve a contraction rate of γ^N at the very lowest. We thus adjust the target contraction rate accordingly, and arrive at the following update:

$$\phi_{k+1} = \phi_k - \varepsilon_k \left(\hat{C}^{(k)}(\alpha(\phi_k)) - \max(\Gamma, \gamma^N) \right). \quad (7)$$

4 Experiments

Having explored the types of trade-offs α -Retrace makes relative to existing off-policy algorithms, we now investigate the performance of these methods in the downstream tasks of evaluation and control described in Section 2.1.

Evaluation. In the left sub-plot of Figure 3, we compare the performance of α -Retrace, n -step uncorrected updates, and n -step importance-weighted updates, for various values of the parameters concerned, at an off-policy evaluation task. In this particular task, the environment is given by a chain MDP (see Appendix Section C.1), the target policy is optimal, and the behaviour is the uniform policy. We plot Q-function L^2 error against number of environment steps; see full details in Appendix Section C.2. Standard error is indicated by the shaded regions.

The best performing methods vary as a function of the number of environment steps experienced. For low numbers of environments steps, the best performing methods n -step uncorrected updates for large n , and α -Retrace for α close to 0. Intuitively, in this regime, a good contraction rate outweighs fixed-point bias. As the number of environment steps increases, the fixed-point bias of the uncorrected methods kicks in, and the optimally-performing α gradually increases from close to 0 to close to 1. Note that typically the extremely high variance of the importance-weighted updates preclude them from attaining any reasonable level of evaluation error.

Control. In the right sub-plot of Figure 3, we compare the performance of a variety of modified policy iteration

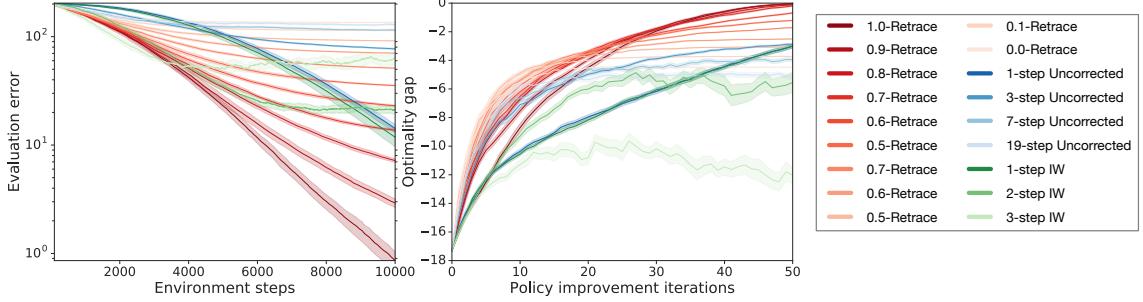


Figure 3: Left: Performance of a variety of off-policy evaluation methods on a small MDP; for further details, see text of Section 4. Right: Performance of a variety of modified policy iteration methods on a small MDP; for further details, see text of Section 4.

methods, each using a different off-policy evaluation method. We use the same MDP as in the evaluation example above, and plot the sub-optimality of the learned policy (measured as difference between expected return under a uniformly random initial state for optimal and learned policies) against the number of policy improvement steps performed. In this experiment, the behaviour policy is fixed as uniform throughout. As with evaluation, we see that initial improvements in policy are strongest with highly-contractive evaluation methods incorporating some fixed-point bias, with less-biased approaches catching up (and ultimately surpassing) when greater amounts of environment interaction are allowed.

4.1 C-trace-R2D2

To test the performance of our methods at scale, we adapted R2D2 [Kapturowski et al., 2019] to replace the original n -step uncorrected update with Retrace and C-trace. For C-trace we targeted the contraction rate given by an n -step uncorrected update, using a discount rate of $\gamma = 0.997$ and $n = 10$. Based on the Pareto efficiency of α -Retrace relative to n -step uncorrected returns exhibited empirically in small-scale MDPs, we conjectured that this should lead to improved performance. The agent was trained on the Atari-57 suite of environments [Bellemare et al., 2013] with the same experimental setup as in [Kapturowski et al., 2019], a description of which we include in Appendix Section E.1. High-level results are displayed in Figure 4, plotting mean human-normalised performance, median human-normalised performance, and mean human-gap (across the 57 games) against wall-clock training time; detailed results are given in Appendix Section F.1.

C-trace-R2D2 attains comparable or superior performance relative to R2D2 and Retrace-R2D2 in all three performance measures. Thus, not only does C-trace-R2D2 match state-of-the-art performance for distributed value-based agents on Atari, it also achieves

the earlier stated goal of bridging the gap between the performance of uncorrected returns and more principled off-policy algorithms in deep reinforcement learning.

4.2 C-trace-DQN

To illustrate the flexibility of C-trace as an off-policy learning algorithm, we also demonstrate its performance within a DQN architecture [Mnih et al., 2015]. We use Double DQN [Van Hasselt et al., 2016] as a baseline, and modify the one-step Q-learning rule to use n -step uncorrected returns, Retrace, and C-trace. As for the R2D2 experiment, we set the C-trace contraction target using $n = 10$, demonstrating the robustness of this C-trace hyperparameter across different architectures. Further, we found the behaviour of C-trace to be generally robust to the choice of n ; see Appendix Section F.2. Full experimental specifications are given in Appendix Section E.2, with detailed results in Appendix Section F.2; a high-level summary is displayed in Figure 4. All sequence-based methods significantly outperform Double DQN, as we would expect. We notice that the performance gap between n -step and Retrace is not as large here as for R2D2 and hypothesize this is due to DQN being much more off-policy than R2D2. As with the R2D2 experiments we see that C-trace-DQN achieves similar learning speed as the targeted n -step update, but with improved final performance. One interpretation of these results is that the improved contraction rate of C-trace allows it to learn significantly faster than Retrace, while the better fixed-point error allows it to find a better long-term solution than n -step uncorrected.

5 Related work

A central observation of this work is that the fixed-point bias can be explicitly traded-off to improve contraction rates. To our understanding, this is the first work to directly study this possibility, and further to draw atten-

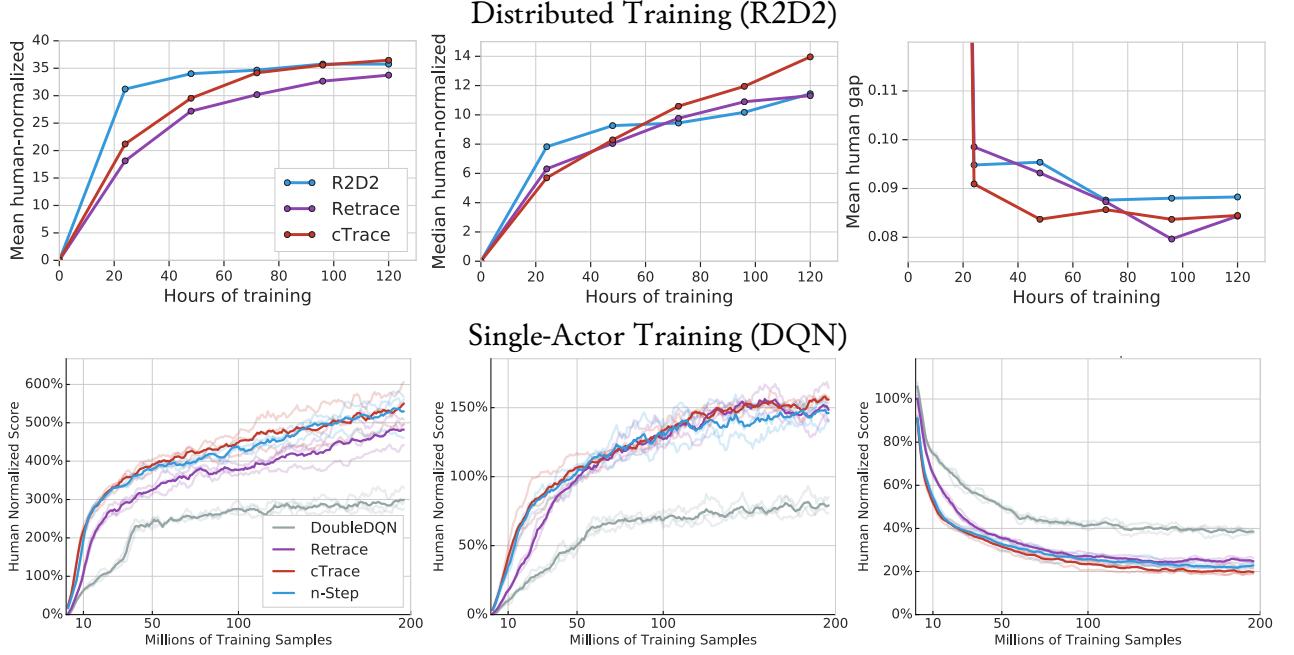


Figure 4: High-level performance of variants of R2D2 (top row) and DQN (bottom row) on the Atari suite of environments. R2D2-based methods are averages of two seeds. DQN-based methods are averages of three seeds. **(Left)** Mean human-normalized score, **(Center)** median human-normalized score, and **(Right)** human gap.

tion to three fundamental quantities to be traded-off in off-policy learning. However, investigating trade-offs in off-policy RL, and in particular parametrising methods to allow a spectrum of algorithms is a long-standing research topic [Sutton and Barto, 2018]. The most closely related methods come from a line of work that consider the bias-variance trade-off due to bootstrapping. In our framework, we understand this as a trade-off between variance and contraction rate, *but without modifying the fixed-point*. The recently introduced $Q(\sigma)$ algorithm uses the σ hyperparameter to mix between importance-weighted n -step SARSA and TreeBackup [De Asis et al., 2018]. In another recent related approach, Shi et al. [2019] uses σ to mix between TreeBackup(λ) and $Q(\lambda)$, although neither of these approaches adaptively set σ based on observed data. We have developed an adaptive method for adjusting α to achieve a desired trace length, and believe an interesting direction for future work would be to develop the adaptive methods described in this paper for use in other families of off-policy learning algorithms.

Conservatively updating policies within control algorithms is a well-established practice; Kakade and Langford [2002] consider a trust-region method for policy improvement, motivated by inexact policy evaluation due to function approximation. In contrast, in this work we consider regularised policy improvement as a means of improving evaluation of future policies, even in

the absence of function approximation. More recently, this approach also led to several advances in policy gradient methods [Schulman et al., 2015, 2017] based on trust regions. Although not the focus of this work, there has been also been much progress on correcting state-visitation distributions [Sutton et al., 2016, Thomas and Brunskill, 2016, Hallak and Mannor, 2017, Liu et al., 2018], another form of off-policy correction important in function approximation, as illustrated in the classic counterexample of Baird [1995].

6 Conclusion

We have highlighted the fundamental role of variance, fixed-point bias, and contraction rate in off-policy learning, and described how existing methods trade off these quantities. With this perspective, we developed novel off-policy learning methods, α -Retrace and C-trace, and incorporated the latter into several deep RL agents, leading to strong empirical performance.

Interesting questions for future work include applying the adaptive ideas underlying C-trace to other families of off-policy algorithms, investigating whether there exist new off-policy learning algorithms in unexplored areas of the space of trade-offs, and developing a deeper understanding of the relationship between these fundamental properties of off-policy learning algorithms and downstream performance on large-scale control tasks.

Acknowledgements

Thanks in particular to Hado van Hasselt for detailed comments and suggestions on an earlier version of this paper, and thanks to Bernardo Avila Pires, Diana Borsa, Steven Kapturowski, Bilal Piot, Tom Schaul, and Yunhao Tang for interesting conversations during the course of this work.

References

- TW Archibald, KIM McKinnon, and LC Thomas. On the generation of Markov decision processes. *Journal of the Operational Research Society*, 46(3):354–361, 1995.
- Leemon Baird. Residual algorithms: Reinforcement learning with function approximation. In *Machine Learning Proceedings*. Elsevier, 1995.
- M. G. Bellemare, Y. Naddaf, J. Veness, and M. Bowling. The arcade learning environment: An evaluation platform for general agents. *Journal of Artificial Intelligence Research*, 47:253–279, jun 2013.
- Dimitri P Bertsekas and John N Tsitsiklis. *Neurodynamic programming*, volume 5. Athena Scientific Belmont, MA, 1996.
- Shalabh Bhatnagar, Richard S Sutton, Mohammad Ghavamzadeh, and Mark Lee. Natural actor–critic algorithms. *Automatica*, 45(11):2471–2482, 2009.
- Kristopher De Asis, J Fernando Hernandez-Garcia, G Zacharias Holland, and Richard S Sutton. Multi-step reinforcement learning: A unifying algorithm. In *AAAI Conference on Artificial Intelligence*, 2018.
- Miroslav Dudík, Dumitru Erhan, John Langford, and Lihong Li. Doubly robust policy evaluation and optimization. *Statistical Science*, 29(4):485–511, 2014.
- Mehrdad Farajtabar, Yinlam Chow, and Mohammad Ghavamzadeh. More robust doubly robust off-policy evaluation. In *International Conference on Machine Learning*, 2018.
- Matthieu Geist and Bruno Scherrer. Off-policy learning with eligibility traces: A survey. *The Journal of Machine Learning Research*, 15(1):289–333, 2014.
- Assaf Hallak and Shie Mannor. Consistent on-line off-policy evaluation. In *International Conference on Machine Learning*, 2017.
- J Fernando Hernandez-Garcia and Richard S Sutton. Understanding multi-step deep reinforcement learning: A systematic study of the DQN target. *arXiv*, 2019.
- 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 *AAAI Conference on Artificial Intelligence*, 2018.
- Sham Kakade and John Langford. Approximately optimal approximate reinforcement learning. In *International Conference on Machine Learning*, 2002.
- Steven Kapturowski, Georg Ostrovski, Will Dabney, John Quan, and Rémi Munos. Recurrent experience replay in distributed reinforcement learning. In *International Conference on Learning Representations*, 2019.
- Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. In *International Conference on Learning Representations*, 2015.
- Qiang Liu, Lihong Li, Ziyang Tang, and Dengyong Zhou. Breaking the curse of horizon: Infinite-horizon off-policy estimation. In *Neural Information Processing Systems*, 2018.
- Ashique Rupam Mahmood, Huizhen Yu, and Richard S Sutton. Multi-step off-policy learning without importance sampling ratios. *arXiv*, 2017.
- Volodymyr Mnih, Koray Kavukcuoglu, David Silver, Andrei A Rusu, Joel Veness, Marc G Bellemare, Alex Graves, Martin Riedmiller, Andreas K Fidjeland, Georg Ostrovski, et al. Human-level control through deep reinforcement learning. *Nature*, 518(7540):529, 2015.
- Rémi Munos, Tom Stepleton, Anna Harutyunyan, and Marc Bellemare. Safe and efficient off-policy reinforcement learning. In *Neural Information Processing Systems*, 2016.
- Bilal Piot, Matthieu Geist, and Olivier Pietquin. Difference of convex functions programming for reinforcement learning. In *Neural Information Processing Systems*, 2014.
- Doina Precup, Rich Sutton, and Satinder Singh. Eligibility traces for off-policy policy evaluation. In *International Conference on Machine Learning*, 2000.
- Herbert Robbins and Sutton Monro. A stochastic approximation method. *Ann. Math. Statist.*, 22(3):400–407, 09 1951.
- John Schulman, Sergey Levine, Pieter Abbeel, Michael I Jordan, and Philipp Moritz. Trust region policy optimization. In *International Conference on Machine Learning*, 2015.
- John Schulman, Filip Wolski, Prafulla Dhariwal, Alec

Radford, and Oleg Klimov. Proximal policy optimization algorithms. *arXiv*, 2017.

Longxiang Shi, Shijian Li, Longbing Cao, Long Yang, and Gang Pan. TBQ(σ): Improving efficiency of trace utilization for off-policy reinforcement learning. In *International Conference on Autonomous Agents and Multiagent Systems*, 2019.

Richard S Sutton and Andrew G Barto. *Reinforcement learning: An introduction*. MIT press, 2018.

Richard S Sutton, A Rupam Mahmood, and Martha White. An emphatic approach to the problem of off-policy temporal-difference learning. *Journal of Machine Learning Research*, 17(1):2603–2631, 2016.

Csaba Szepesvári. Algorithms for reinforcement learning. *Synthesis lectures on artificial intelligence and machine learning*, 4(1):1–103, 2010.

Philip Thomas and Emma Brunskill. Data-efficient off-policy policy evaluation for reinforcement learning. In *International Conference on Machine Learning*, 2016.

Hado Van Hasselt, Arthur Guez, and David Silver. Deep reinforcement learning with double Q-learning. In *AAAI Conference on Artificial Intelligence*, 2016.

Ziyu Wang, Tom Schaul, Matteo Hessel, Hado van Hasselt, Marc Lanctot, and Nando de Freitas. Dueling network architectures for deep reinforcement learning. In *International Conference on Machine Learning*, 2016.

Appendices: Adaptive Trade-Offs in Off-Policy Learning

A Proofs

Proposition 2.1. Consider the task of evaluation of a policy π under behaviour μ , and consider an update rule \hat{T} which stochastically approximates the application of an operator \tilde{T} , with contraction rate Γ and fixed point \tilde{Q} , to an initial estimate Q . Then we have the following decomposition:

$$\begin{aligned} \mathbb{E} [\|\hat{T}Q - Q^\pi\|_\infty] &\leq \\ \underbrace{\mathbb{E} [\|\hat{T}Q - \tilde{T}Q\|_2^2]}_{(\text{Root}) \text{ variance}}^{1/2} + \underbrace{\Gamma \|Q - \tilde{Q}\|_\infty}_{\text{Contraction}} + \underbrace{\|\tilde{Q} - Q^\pi\|_2}_{\text{Fixed-point bias}}. \end{aligned}$$

Proof. Note that by the triangle inequality:

$$\mathbb{E} [\|\hat{T}Q - Q^\pi\|_\infty] \leq \mathbb{E} [\|\hat{T}Q - \tilde{T}Q\|_\infty] + \|\tilde{T}Q - \tilde{Q}\|_\infty + \|\tilde{Q} - Q^\pi\|_\infty.$$

Now, observing $\|\tilde{T}Q - \tilde{Q}\|_\infty = \|\tilde{T}Q - \tilde{T}\tilde{Q}\|_\infty \leq \Gamma \|Q - \tilde{Q}\|_\infty$ yields the second term on the right-hand side of the stated bound. Using the inequality $\|\cdot\|_\infty \leq \|\cdot\|_2$ and Jensen's inequality yields the remaining terms. \square

Theorem 2.2. Consider an update rule \hat{T} with corresponding operator T , and consider the collection $\mathcal{M} = \mathcal{M}(\mathcal{X}, \mathcal{A}, P, \gamma, R_{\max})$ of MDPs with common state space, action space, transition kernel, and discount factor (but varying rewards, with maximum immediate reward bounded in absolute value by R_{\max}). Fix a target policy π , and a random variable Z , the set of transitions used by the operator \hat{T} ; these could be transitions encountered in a trajectory following the behaviour μ , or i.i.d. samples from the discounted state-action visitation distribution under μ . We denote the mismatch between π and Z at level $\delta \in (0, 1)$ by

$$\begin{aligned} D(Z, \pi, \delta) &\stackrel{\text{def}}{=} \max \{d_{(x,a),\pi}(\Omega) \mid \Omega \subseteq \mathcal{X} \times \mathcal{A} \text{ s.t.} \\ &\quad \mathbb{P}(Z \cap \Omega \neq \emptyset) \leq \delta, (x, a) \in \mathcal{X} \times \mathcal{A}\}, \end{aligned}$$

where $d_{(x,a),\pi}$ is the discounted state-action visitation distribution for trajectories initialised with (x, a) , following π thereafter. Denoting the variance, contraction rate, and fixed-point bias of \hat{T} for a particular MDP $M \in \mathcal{M}$ by $\mathbb{V}(M)$, $\Gamma(M)$ and $B(M)$ respectively, we have

$$\begin{aligned} \sup_{M \in \mathcal{M}} \left[\sqrt{\mathbb{V}(M)} + \frac{2R_{\max}}{1-\gamma} \Gamma(M) + B(M) \right] &\geq \\ \sup_{\delta \in (0,1)} (1-\delta) D(Z, \pi, \delta) R_{\max} / (1-\gamma). \end{aligned}$$

Proof. The high-level approach to the proof is to exhibit two MDPs $M_0, M_1 \in \mathcal{M}$ which with high probability under the data used by \hat{T} , cannot be distinguished. This yields a high probability lower bound on the evaluation error that the operator \hat{T} achieves on the two MDPs. This in turn implies that the mean-squared error quantity of Proposition 2.1 cannot be uniformly low across M_0 and M_1 , and this yields a lower bound for the quantity on the right-hand side of the bound appearing in Proposition 2.1, as required.

Using the notation introduced in the statement of the theorem, for a given $\delta \in (0, 1)$, let $(x^*, a^*) \in \mathcal{X} \times \mathcal{A}$, $\Omega^* \subseteq \mathcal{X}$ be quantities achieving the maximum in the definition of $D(Z, \pi, \delta)$. Thus, with probability at least $1 - \delta$, none of the state-action pairs used by the algorithm \hat{T} are contained in Ω^* .

Now define two MDPs M_0, M_1 with common state space \mathcal{X} , action space \mathcal{A} , transition kernel P , and discount factor γ , with reward functions $r_0, r_1 : \mathcal{X} \times \mathcal{A} \rightarrow \mathbb{R}$ defined by

$$r_i(x, a) = \begin{cases} 0 & \text{if } x \notin \Omega^* \\ (-1)^i R_{\max} & \text{if } x \in \Omega^*. \end{cases}$$

Now, the Q-functions associated with these two MDPs can be calculated by $(I - \gamma P^\pi)^{-1} r_0$ and $(I - \gamma P^\pi)^{-1} r_1$, and so we can read off their difference in the (x^*, a^*) coordinate as

$$\sum_{(x,a) \in \Omega^*} \frac{1}{1-\gamma} d_{(x^*, a^*), \pi}(x, a) 2R_{\max} = \frac{2D(Z, \pi, \delta) R_{\max}}{1-\gamma}.$$

Thus, with probability $1 - \delta$, the algorithm must make an error of at least $D(Z, \pi, \delta) R_{\max} / (1 - \gamma)$ on one of the MDPs M_0 and M_1 , as measured the L^∞ norm. This implies that the expected L^∞ error appearing on the left-hand side of the bound in Proposition 2.1 is at least $(1 - \delta) D(Z, \pi, \delta) R_{\max} / (1 - \gamma)$ for one of the MDPs M_0 and M_1 . Thus, the statement of the theorem follows by taking a supremum over $\delta \in (0, 1)$. \square

Proposition 3.2. The operator associated with the α -Retrace evaluation update for evaluating π given behaviour μ has a state-action-dependent contraction rate of

$$C(\alpha|x, a) \stackrel{\text{def}}{=} 1 - (1 - \gamma) \times \mathbb{E}_\mu \left[\sum_{t=0}^{\infty} \gamma^t \prod_{s=1}^t ((1 - \alpha) + \alpha \bar{\rho}_s) \middle| (X_0, A_0) = (x, a) \right], \quad (5)$$

for each $(x, a) \in \mathcal{X} \times \mathcal{A}$. Viewed as a function of $\alpha \in [0, 1]$, this contraction rate is continuous, monotonically increasing, with minimal value 0, and maximal value no greater than γ . Further, the contraction rate is *strictly* monotonic iff π and μ are (x, a) -distinguishable under μ .

Proof. The α -Retrace operator for evaluation of π given behaviour μ corresponds to the standard Retrace operator for evaluation of $\pi^\alpha = \alpha\pi + (1 - \alpha)\mu$ given behaviour μ . Thus, from the analysis of Munos et al. [2016], the contraction rate of the α -Retrace operator specific to a particular state-action pair $(x, a) \in \mathcal{X} \times \mathcal{A}$ may be immediately written down as

$$\begin{aligned} & 1 - (1 - \gamma) \mathbb{E}_\mu \left[\sum_{t \geq 0} \gamma^t \prod_{s=1}^t \min \left(1, \frac{\alpha\pi(A_t|X_t) + (1 - \alpha)\mu(A_t|X_t)}{\mu(A_t|X_t)} \right) \middle| X_0 = x, A_0 = a \right] \\ &= 1 - (1 - \gamma) \mathbb{E}_\mu \left[\sum_{t \geq 0} \gamma^t \prod_{s=1}^t \left((1 - \alpha) + \alpha \min \left(1, \frac{\pi(A_t|X_t)}{\mu(A_t|X_t)} \right) \right) \middle| X_0 = x, A_0 = a \right]. \end{aligned}$$

To see that this is a continuous function of $\alpha \in [0, 1]$, we note that the integrand of the expectation above is clearly a continuous function of α , and is uniformly dominated by the constant function equal to $(1 - \gamma)^{-1}$. By the dominated convergence theorem, continuity of the above expression follows. Since the integrand is non-negative and bounded above by $(1 - \gamma)^{-1}$, the contraction rate must lie in the interval $[0, \gamma]$ for all $\alpha \in [0, 1]$. For monotonicity, we show that each term

$$\mathbb{E}_\mu \left[\prod_{s=1}^t \left((1 - \alpha) + \alpha \min \left(1, \frac{\pi(A_t|X_t)}{\mu(A_t|X_t)} \right) \right) \middle| X_0 = x, A_0 = a \right] \quad (8)$$

is monotonic decreasing in α , meaning that the contraction rate is monotonic increasing in α . To this end, observe that the integrand of the expectation above almost-surely takes the form $\prod_{s=1}^t (1 - z_s \alpha)$ for some coefficients $z_s \in [0, 1]$. The derivative with respect to α of this expression is $\sum_{s=1}^t -z_s \prod_{s' \neq s}^t (1 - z_{s'} \alpha)$, which is non-positive for $\alpha \in [0, 1]$. It is again straightforward to apply the dominated convergence theorem to this derivative to obtain that the derivative of Expression (8) is non-positive for all $\alpha \in [0, 1]$, and we thus obtain monotonicity as required. Finally, for strict monotonicity, note that if π and μ are *not* distinguishable under μ , then all truncated importance weights in the expressions above are equal to 1 almost-surely under the distribution over states visited when following μ . Hence, the contraction rate is in fact constant as a function of α , and we therefore do not have strict monotonicity. On the other hand, if π and μ are (x, a) -distinguishable under μ , then there exists a $t \in \mathbb{N}$ such that in the integrand of the expectation in Expression (8), in one of the terms constituting the product, the coefficient of α is less than 0 with positive probability. Thus, the integrand is strictly monotonic with positive probability, and hence Expression (8) itself is strictly monotonic, proving strict monotonicity of the contraction rate, as required. \square

Proposition 3.3. Consider a target policy π , let μ be the behavioural policy, and assume that there is a unique greedy action $a^*(x) \in \mathcal{A}$ with respect to Q^π at each state x for each $x \in \mathcal{X}$. Then there exists a value of $\alpha \in (0, 1)$ such that the greedy policy with respect to the fixed point of α -Retrace coincides with the greedy policy with respect to Q^π , and the contraction rate for this α -Retrace is no greater than that for 1-Retrace. Further, if π and μ are (x, a) -distinguishable under μ for all $(x, a) \in \mathcal{X} \times \mathcal{A}$, then the contraction rate of α -Retrace is strictly lower than that of 1-Retrace.

Proof. That the greedy policies coincide follows as a consequence of the continuity of Q^ν with respect to the policy ν and the positivity of the minimum action gap $\Delta = \inf_{x \in \mathcal{X}, a \neq a^*(x)} (Q^\pi(x, a^*(x)) - Q^\pi(x, a))$; α may be selected so that e.g. $\|Q^{\pi^\alpha} - Q^\pi\|_\infty \leq \Delta/2$. The contraction result follows from the monotonicity property derived in Proposition 3.2. \square

Proposition 3.4. Let $(x_t^{(k)}, a_t^{(k)}, r_t^{(k)})_{t=0}^\infty$ be an i.i.d. sequence of trajectories following μ , with initial state-action distribution given by ν . Let Γ be a target contraction rate such that $C_\nu(1) \geq \Gamma$. Let the stepsizes $(\varepsilon_k)_{k=0}^\infty$ satisfy the usual Robbins-Monro conditions $\sum_{k=0}^\infty \varepsilon_k = \infty$, $\sum_{k=0}^\infty \varepsilon_k^2 < \infty$. Then for any initial value ϕ_0 following the updates in (6), we have $\phi_k \rightarrow \phi^*$ in probability, where $\phi^* \in \mathbb{R}$ is the unique value such that $C_\nu(\alpha(\phi^*)) = \Gamma$.

Proof. The proof follows from an application of standard stochastic approximation theory to the solution of the root-finding problem $C_\nu(\alpha(\phi)) = \Gamma$. Firstly, by Proposition 3.2, the function $\phi \mapsto C_\nu(\alpha(\phi))$ is continuous and monotonic on \mathbb{R} . By the assumption that $C_\nu(1) \geq \Gamma$, it follows that $\phi \mapsto C_\nu(\alpha(\phi))$ is strictly monotonic, and moreover by inspecting the proof of Proposition 3.2, has positive derivative everywhere. By the intermediate value theorem, there exists a unique value $\phi^* \in \mathbb{R}$ such that $C_\nu(\alpha(\phi^*)) = \Gamma$.

Now note that for each $\phi \in \mathbb{R}$, the random variables $\hat{C}^{(k)}(\alpha(\phi))$ are i.i.d. unbiased, bounded estimators of $C_\nu(\alpha(\phi))$. Thus, the scheme (6) is a standard stochastic approximation scheme for the the root of a monotonic function, and the conditions of Theorem 2 of [Robbins and Monro, 1951] are satisfied, enabling us to conclude that $\phi_t \rightarrow \phi^*$ in probability, as required. \square

Theorem 3.5. Assume the same conditions as Proposition 3.4, and additionally that: (i) trajectory lengths have finite second moment; (ii) immediate rewards are bounded. Let $(\phi_k)_{k=0}^\infty$ be defined as in Equation (6) and $(Q_k)_{k=0}^\infty$ be a sequence of Q-functions, with Q_{k+1} obtained from applying Retrace updates targeting $\alpha(\phi_k)\pi + (1 - \alpha(\phi_k))\mu$ to Q_k with trajectory $k + 1$, using stepsize ε_k . Then we have $\alpha(\phi_k) \rightarrow \alpha(\phi^*) =: \alpha^*$ and $Q_k \rightarrow Q^{\alpha^*\pi+(1-\alpha^*)\mu}$ almost surely.

Proof. Convergence in probability of $\alpha_k := \alpha(\phi_k)$ to α^* has been shown in Proposition 3.4; it is straightforward to upgrade this to almost-sure convergence using standard stochastic approximation theory. The intuition for the remainder of the proof is that when α_k is close to α^* , the C-trace updates are close to those of standard Retrace targeting the policy $\alpha^*\pi + (1 - \alpha^*)\mu$, which are known to converge under the conditions of the theorem. This is made rigorous by decomposing the update on the Q-function from the $(k + 1)^{\text{th}}$ trajectory as

$$Q_{k+1} = \underbrace{(\mathbf{1} - \tilde{\varepsilon}_k) \odot Q_k + \tilde{\varepsilon}_k \odot \mathcal{R}^{\alpha^*} Q_k}_{\text{Desired update}} + \underbrace{(Q_{k+1} - (\mathbf{1} - \tilde{\varepsilon}_k) \odot Q_k - \tilde{\varepsilon}_k \odot \mathcal{R}^{\alpha_k} Q_k)}_{\text{Martingale noise}} + \underbrace{\tilde{\varepsilon}_k \odot (\mathcal{R}^{\alpha_k} Q_k - \mathcal{R}^{\alpha^*} Q_k)}_{\text{Perturbation}},$$

where \mathcal{R}^α denotes the Retrace operator targeting $\alpha\pi + (1 - \alpha)\mu$, and with $\tilde{\varepsilon}_k(x, a) = \varepsilon_k \mathbb{E}[\sum_t \mathbb{1}_{(x_t, a_t)=(x, a)} | (x_0, a_0) = (x, a)]$, and \odot the Hadamard product and $\mathbf{1}$ the vector of 1's. It is then possible to appeal to Proposition 4.5 of Bertsekas and Tsitsiklis [1996] that $Q_k \rightarrow Q^{\alpha^*\pi+(1-\alpha^*)\mu}$ almost surely, using the assumptions of theorem. \square

B Additional results

B.1 Operators for time-inhomogeneous policies

In this section, we provide a result which rigorously proves the connection between the n -step uncorrected target and the time-inhomogeneous policy mentioned in Section 2.

Proposition B.1. The n -step uncorrected update corresponding to the target

$$\sum_{s=0}^{n-1} \gamma^s r_s + \gamma^n \mathbb{E}_{A \sim \pi(\cdot|x_n)} [\hat{Q}(x_n, A)],$$

with the trajectory generated under μ , is a stochastic approximation to the operator $(T^\mu)^{n-1} T^\pi$, with fixed point given by the Q -function for the time-inhomogeneous policy which follows π at timesteps t satisfying $t \equiv n-1 \pmod{n}$, and μ otherwise.

Proof. We begin by taking the expectation of the update target conditional on the initial state-action pair, and showing that it is equal to $((T^\mu)^{n-1} T^\pi \hat{Q})(x_0, a_0)$. We proceed by induction. In the case $n = 1$, the expectation of the update is given by

$$\begin{aligned} & \mathbb{E}_\mu \left[R(X_0, A_0) + \gamma \mathbb{E}_{A \sim \pi(\cdot|X_1)} [\hat{Q}(X_1, A)] \middle| X_0 = x_0, A_0 = a_0 \right] \\ &= r(x_0, a_0) + \sum_{x' \in \mathcal{X}} P(x'|x, a) \gamma \sum_{a' \in \mathcal{A}} \pi(a'|x') \hat{Q}(x', a') \\ &= (T^\pi \hat{Q})(x_0, a_0), \end{aligned}$$

as required. For the inductive step, we assume the result holds for some $n \geq 1$. Now observe that by conditioning on (X_1, A_1) , we have

$$\begin{aligned} & \mathbb{E}_\mu \left[\sum_{s=0}^n \gamma^s R(X_s, A_s) + \gamma^{n+1} \mathbb{E}_{A \sim \pi(\cdot|X_n)} [\hat{Q}(X_{n+1}, A)] \middle| X_0 = x_0, A_0 = a_0 \right] \\ &= r(x_0, a_0) + \gamma \sum_{x_1 \in \mathcal{X}} P(x_1|x_0, a_0) \sum_{a_1 \in \mathcal{A}} \mu(a_1|x_1) \times \\ & \quad \mathbb{E} \left[\sum_{s=1}^n \gamma^{s-1} R(X_s, A_s) + \gamma^{n+1} \mathbb{E}_{A \sim \pi(\cdot|X_n)} [\hat{Q}(X_{n+1}, A)] \middle| X_1 = x_1, A_1 = a_1 \right] \\ &\stackrel{(a)}{=} r(x_0, a_0) + \gamma \sum_{x_1 \in \mathcal{X}} P(x_1|x_0, a_0) \sum_{a_1 \in \mathcal{A}} \mu(a_1|x_1) ((T^\mu)^{n-1} T^\pi \hat{Q})(x_1, a_1) \\ &= (T^\mu (T^\mu)^{n-1} T^\pi \hat{Q})(x_0, a_0) \\ &= ((T^\mu)^n T^\pi \hat{Q})(x_0, a_0), \end{aligned}$$

as required, with (a) following from the induction hypothesis. Finally, for the interpretation of the fixed point of $(T^\mu)^{n-1} T^\pi$, observe that the time-inhomogeneous policy described in the statement of the proposition, which we denote $\pi \mu^{n-1}$ follows a stream of Markovian policies with period n , so it is possible to write down an n -step Bellman equation for its Q -function $Q^{\pi \mu^{n-1}}$. Doing so yields

$$\begin{aligned} Q^{\pi \mu^{n-1}}(x, a) &= \mathbb{E}_{\substack{A_{1:n-1} \sim \mu(\cdot|X_{1:n-1}) \\ A_n \sim \pi(\cdot|X_n)}} \left[\sum_{s=0}^{n-1} \gamma^s R(X_s, A_s) + \gamma^n Q^{\pi \mu^{n-1}}(X_n, A_n) \middle| X_0 = x, A_0 = a \right] \\ &= \mathbb{E}_\mu \left[\sum_{s=0}^{n-1} \gamma^s R(X_s, A_s) + \gamma^n \mathbb{E}_{A_n \sim \pi(\cdot|X_n)} [Q^{\pi \mu^{n-1}}(X_n, A_n)] \middle| X_0 = x, A_0 = a \right]. \end{aligned}$$

We recognise the right-hand side as the operator $(T^\mu)^{n-1} T^\pi$, and thus $Q^{\pi \mu^{n-1}}$ is the fixed point of this operator. \square

B.2 Further decompositions of evaluation error

In addition to the decomposition given in Proposition 2.1, there are decompositions of evaluation error based on other norms that may be of interest. We state one such decomposition below, and also note that there is also scope to use different norms to define the fundamental traded-off quantities, such as using the L^∞ norm to define an alternative notion of fixed-point bias, that lead to further decompositions.

Proposition B.2. Consider the task of evaluation of a policy π under behaviour μ , and consider an update rule \hat{T} which stochastically approximates the application of an operator T , with contraction rate Γ and fixed point \tilde{Q} , to an initial estimate Q . Then we have the following decomposition:

$$\mathbb{E} [\|\hat{T}Q - Q^\pi\|_2^2] \leq 3 \left[\underbrace{\mathbb{E} [\|\hat{T}Q - TQ\|_2^2]}_{\text{Variance}} + \underbrace{\Gamma^2 |\mathcal{X}| |\mathcal{A}| \|Q - \tilde{Q}\|_\infty^2}_{(\text{Squared}) \text{ contraction}} + \underbrace{\|\tilde{Q} - Q^\pi\|_2^2}_{(\text{Squared}) \text{ fixed-point bias}} \right].$$

Proof. The inequality is obtained in a manner analogous to that of Proposition 2.1. First, a Cauchy-Schwarz-style argument yields

$$\mathbb{E} [\|\hat{T}Q - Q^\pi\|_2^2] \leq 3 \left[\mathbb{E} [\|\hat{T}Q - \tilde{T}Q\|_2^2] + \|\tilde{T}Q - \tilde{Q}\|_2^2 + \|\tilde{Q} - Q^\pi\|_2^2 \right].$$

Then, the inequality $\|\cdot\|_2 \leq |\mathcal{X}| |\mathcal{A}| \|\cdot\|_\infty$ is applied, together with the definition of T as a contraction mapping under $\|\cdot\|_\infty$ with contraction modulus Γ , to yield the statement. \square

C Experimental details

C.1 Environments

Dirichlet-Uniform random MDPs. These random MDPs are specified by two parameters: the number of states, n_s , and the number of actions, n_a . Transition distributions $P(\cdot|x, a)$ are sampled i.i.d. from a Dirichlet(1, ..., 1) distribution for each $\mathcal{X} \times \mathcal{A}$. Each immediate reward distribution is given by a Dirac delta, with locations drawn i.i.d. from the Uniform([-1, 1]) distribution.

Garnet MDPs. Garnet MDPs [Archibald et al., 1995, Piot et al., 2014, Bhatnagar et al., 2009, Geist and Scherrer, 2014] are drawn from a distribution specified by three numbers: the number of states, n_s , the number of actions, n_a , and the *branching factor*, n_b . Each transition distribution $P(\cdot|x, a)$ is given by $n_b^{-1} \sum_{i=1}^{n_b} \delta_{z_i(x, a)}$, where $z_{1:n_b}(x, a)$ are drawn uniformly without replacement from the set of states of the MDP, independently for each state-action pair $(x, a) \in \mathcal{X} \times \mathcal{A}$. $\lfloor n_s / 10 \rfloor$ states are selected uniformly without replacement, such that any transition out of these states yields a reward of 1, whilst all other transitions in the MDP yield a reward of 0.

Chain MDP. Our chain MDP is specified by a number of states n_s , identified with the set $\{1, \dots, n_s\}$. State n_s is terminal. Two actions, `left` and `right`, are available at each state, which deterministically move the agent into the corresponding state (taking action `left` in state 1 causes the agent to remain in state 1). Every transition caused by the action `right` incurs a reward of -1 , unless the transition is into state n_s , in which case a reward of 50 is received.

C.2 Additional details for plots appearing in the main paper

Figure 1. We use a Dirichlet-Uniform random MDP (see Section C.1) with 5 states and 3 actions. The target π and behaviour μ policies were sampled independently, so that each distribution $\pi(\cdot|x)$ and $\mu(\cdot|x)$ are independent draws from the Dirichlet(1, ..., 1) distribution. We use a discount rate of 0.9, and a uniform initial state distribution. The variance variable is estimated by simulating 5000 trajectories of length 100, from an initial Q-function estimate set to 0.

Figure 3. In both tasks, the environment is the chain described in Section C.1 with $n_s = 20$. In both tasks, all learning algorithms use a learning rate of 0.1, and the discount factor is set to 0.9 throughout. In the control task, policy improvement is interleaved with 100 steps of environment experience, which are used by the relevant evaluation algorithm. All Retrace-derived methods use $\lambda = 1$. In both evaluation and control tasks, the experiments were repeated 200 times to estimate the standard error by bootstrapping, which is indicated in the plots by the shaded regions.

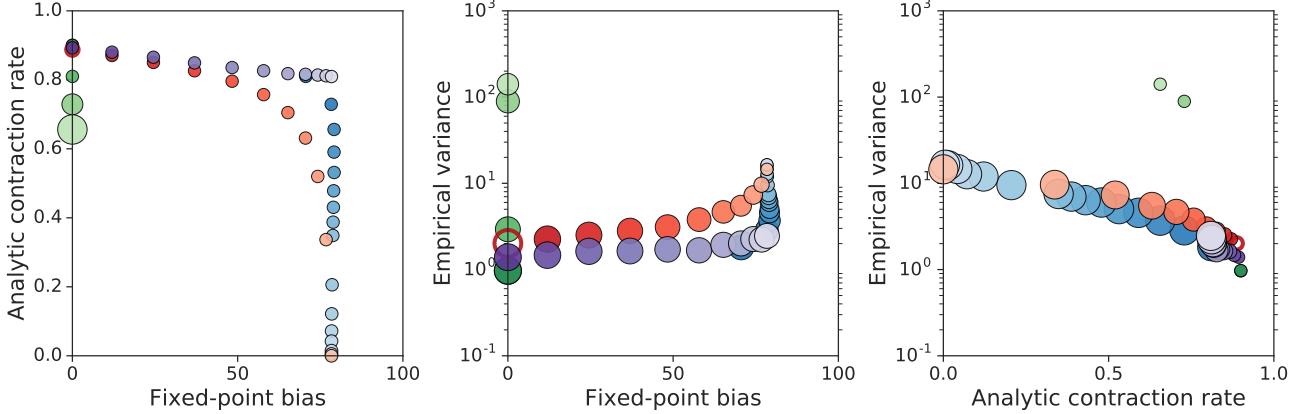


Figure 5: Trade-offs made by n -step uncorrected methods ($n = 1$ (light blue) through to $n = 50$ (dark blue)), n -step importance-weighted methods ($n = 1$ (dark green) through to $n = 4$ (dark green), α -Retrace ($\alpha = 1$ (dark red) through to $\alpha = 0$ (light red)), and α -TreeBackup ($\alpha = 1$ (dark purple) through to $\alpha = 0$ (light purple)). Results are shown for the chain environment, and evaluation of a Dirichlet($1, \dots, 1$) policy under behaviour generated by an independently sampled Dirichlet($1, \dots, 1$) policy.

D Further experimental results

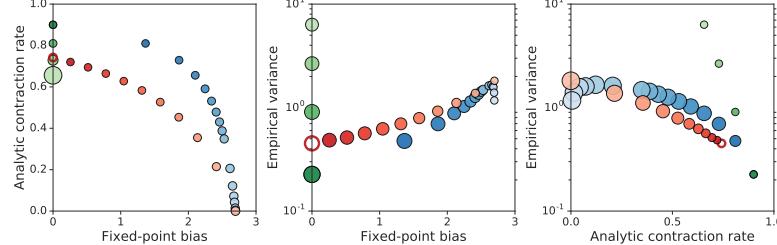
D.1 Further trade-off plots

In this section, we give several further examples of trade-offs made by off-policy algorithms. We begin by examining the trade-offs made by TreeBackup, for which the update target (for a target policy π given a trajectory generated according to a behaviour policy μ) is stated below for completeness.

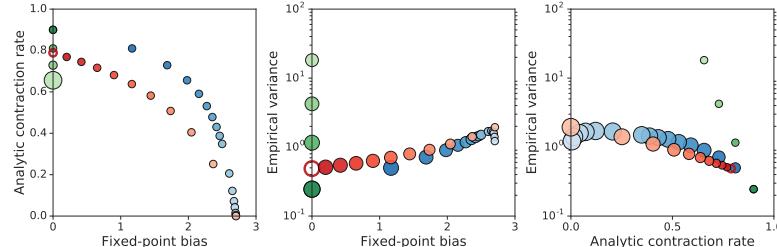
$$\hat{Q}(x_0, a_0) + \sum_{s \geq 0} \gamma^s \prod_{u=1}^s \pi(a_u | x_u) \left(r_s + \gamma \mathbb{E}_{A \sim \pi(\cdot | x_{s+1})} [\hat{Q}(x_{s+1}, A)] - \hat{Q}(x_s, a_s) \right).$$

We show that mixing in a proportion $1 - \alpha$ of the behaviour policy into the target in TreeBackup (which we dub α -TreeBackup) leads to fundamentally different trade-off behaviour than in α -Retrace; see Figure 5. As can be seen in the plot, mixing in the behaviour policy leads to limited improvements in contraction rate relative to the trade-off achieved by α -Retrace, whilst incurring significant fixed-point bias.

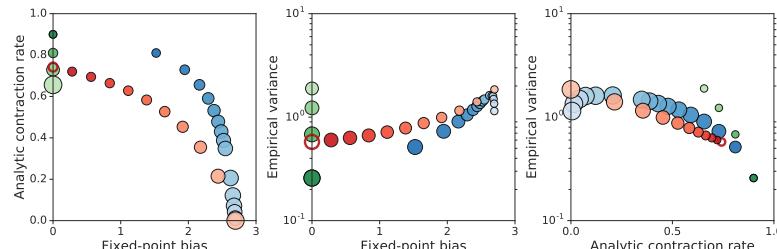
We next demonstrate the robustness of the behaviour exhibited in Figure 1 in a variety of environments, and with a variety of target/behaviour policy pairings. As in Figure 1, α -Retrace is illustrated in red, with dark red corresponding to $\alpha = 1$ through to $\alpha = 0$ in light red. n -step uncorrected methods are illustrated in blue, ranging from $n = 1$ (dark blue) through to $n = 50$ (light blue). n -step importance-weighted methods are illustrated in green, ranging from $n = 1$ (dark green) through to $n = 4$ (light green). Results are given for a Dirichlet-Uniform random MDP (Figure 6), a random garnet MDP (Figure 7), and the chain MDP described in Section C.1 (Figure 8). In all cases, we use a discount rate $\gamma = 0.9$, a learning rate for each algorithm of 0.1, and the variance variable is estimated from 5000 i.i.d. trajectories of length 100. All Retrace methods use $\lambda = 1$ (as presented in the main paper).



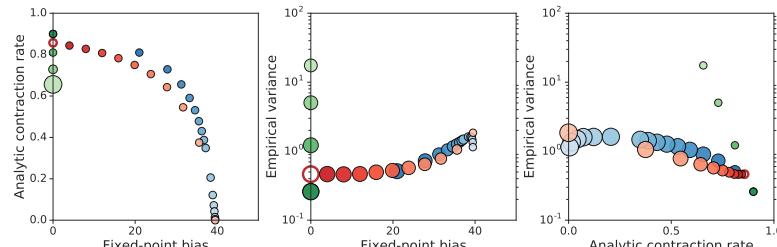
(a) Target policy: uniform. Behaviour policy: Dirichlet(1, ..., 1) random.



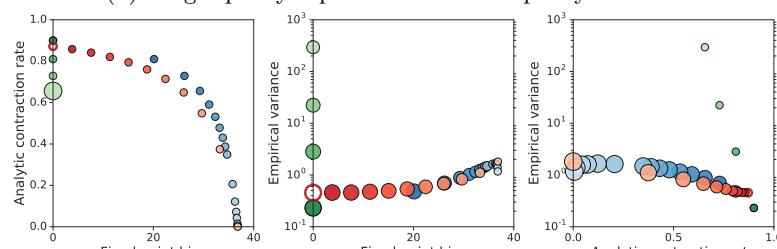
(b) Target policy: Dirichlet(1, ..., 1) random. Behaviour policy: Independent Dirichlet(1, ..., 1) random.



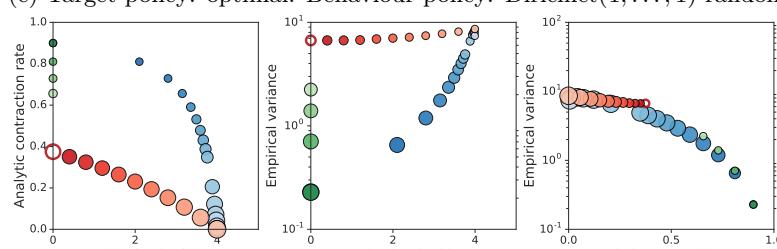
(c) Target policy: Dirichlet(1, ..., 1) random. Behaviour policy: uniform.



(d) Target policy: optimal. Behaviour policy: uniform.

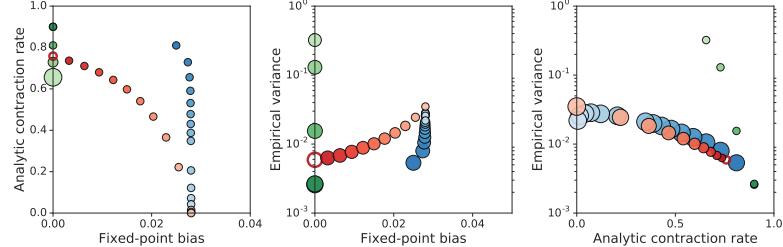


(e) Target policy: optimal. Behaviour policy: Dirichlet(1, ..., 1) random.

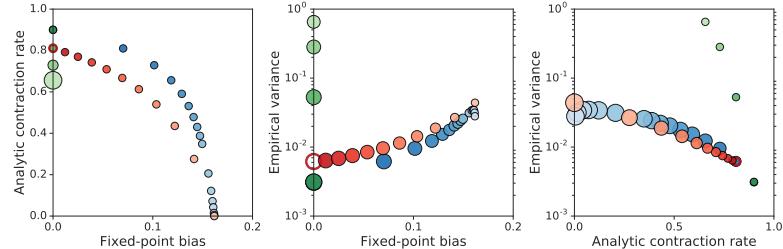


(f) Target policy: optimal. Behaviour policy: optimal, with uniform exploration at probability 0.1.

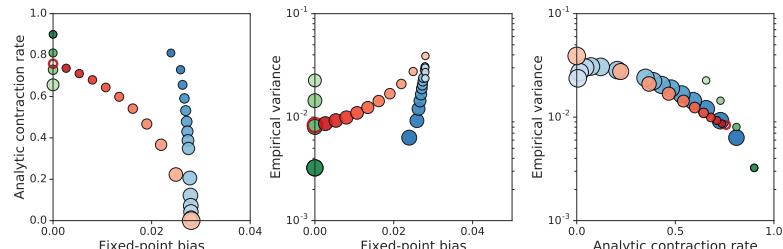
Figure 6: Trade-off plots for a Dirichlet-Uniform random MDP with 20 states and 3 actions, with a variety of target policy/behaviour policy pairings.



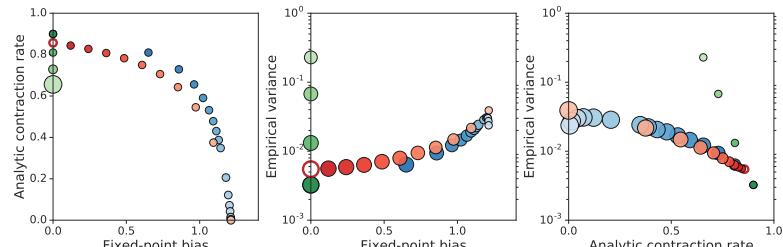
(a) Target policy: uniform. Behaviour policy: Dirichlet($1, \dots, 1$) random.



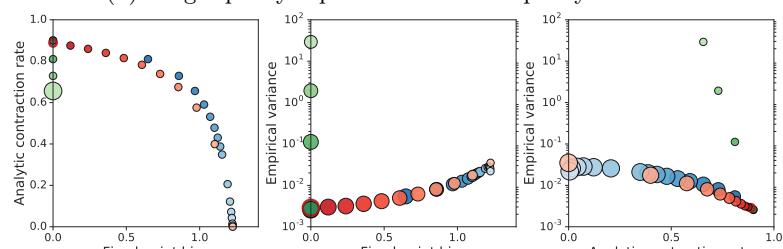
(b) Target policy: Dirichlet($1, \dots, 1$) random. Behaviour policy: Independent Dirichlet($1, \dots, 1$) random.



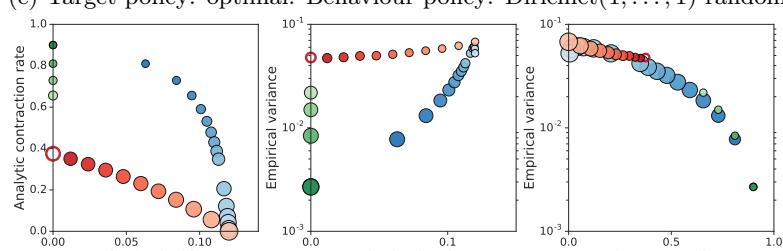
(c) Target policy: Dirichlet($1, \dots, 1$) random. Behaviour policy: uniform.



(d) Target policy: optimal. Behaviour policy: uniform.

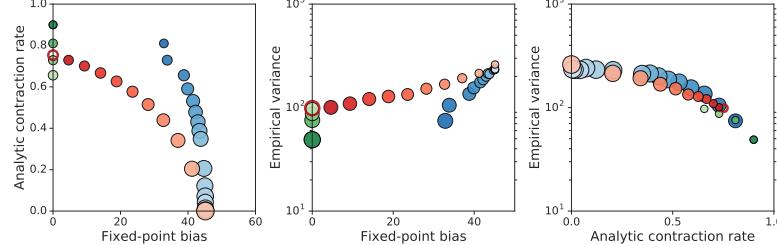


(e) Target policy: optimal. Behaviour policy: Dirichlet($1, \dots, 1$) random.

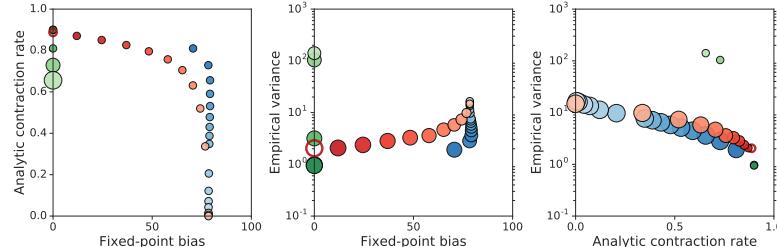


(f) Target policy: optimal. Behaviour policy: optimal, with uniform exploration at probability 0.1.

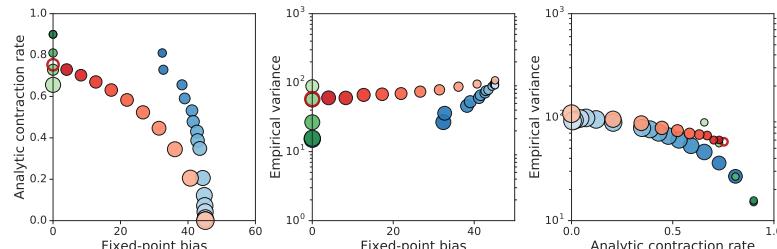
Figure 7: Trade-off plots for a garnet random MDP with 20 states and 3 actions, with a variety of target policy/behaviour policy pairings.



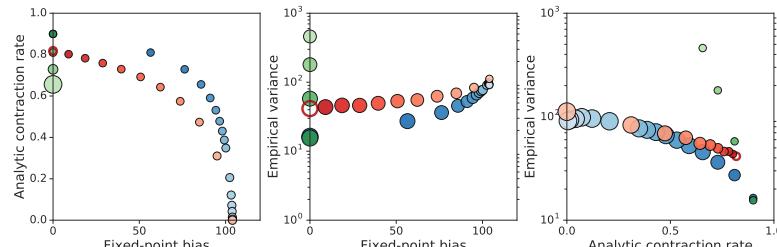
(a) Target policy: uniform. Behaviour policy: Dirichlet(1, ..., 1) random.



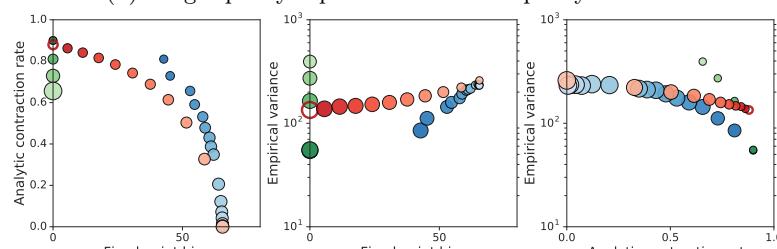
(b) Target policy: Dirichlet(1, ..., 1) random. Behaviour policy: Independent Dirichlet(1, ..., 1) random.



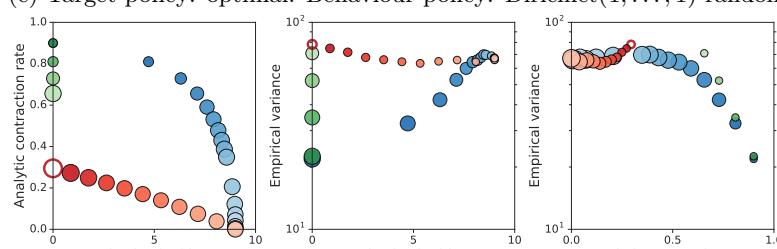
(c) Target policy: Dirichlet(1, ..., 1) random. Behaviour policy: uniform.



(d) Target policy: optimal. Behaviour policy: uniform.



(e) Target policy: optimal. Behaviour policy: Dirichlet(1, ..., 1) random.



(f) Target policy: optimal. Behaviour policy: optimal, with uniform exploration at probability 0.1.

Figure 8: Trade-off plots for the chain MDP described in Section C.1, with 20 states, with a variety of target policy/behaviour policy pairings.

E Large-scale experiment details

Episodes are limited to 30 minutes (108,000 environment frames). When reporting numeric scores, as opposed to learning curves, we give final agent performance as undiscounted episodic returns. The computing architectures used to run the two agents correspond precisely to the descriptions given in Mnih et al. [2015], Kapturowski et al. [2019].

Mini-batches are drawn from an experience replay buffer as described in the baseline agent papers [Mnih et al., 2015, Kapturowski et al., 2019]. For Retrace and C-trace, the n -step loss function is modified to use the Retrace update for the, possibly modified, target policy. GPU training was performed on an NVIDIA Tesla V100.

Only for R2D2 experiments, all agents (including Retrace-based algorithms) use the invertible value function rescaling of R2D2. Finally, for C-trace, the target policy is given by

$$\hat{\pi} := (1 - \alpha)\pi + \alpha\mu,$$

where π is the greedy policy on the current action-values and μ is the ϵ -greedy policy followed by the actor generating the current trajectory. The value of α is adapted with each mini-batch using Robbins-Monro updates with truncated trajectory targets, as described in Section (3.2). The average observed contraction rate of Retrace over the mini-batch is calculated from the Retrace weights (see Equation (5)),

$$\hat{C}(\alpha) = 1 - (1 - \gamma) \sum_{t=0}^N \gamma^t \prod_{s=1}^t \left((1 - \alpha) + \alpha \min \left(1, \frac{\pi(a_s|x_s)}{\mu(a_s|x_s)} \right) \right).$$

For simplicity we restate the Robbins-Monro update as a loss, scale up by $1000/(1 - \gamma)$ (to counter-act the small learning rate from Adam) and add it to the primary loss. We use a value of $\lambda = 1.0$ for all Retrace and C-trace large-scale experiments. We considered $\lambda = 0.97$, in keeping with published work, but found larger values to perform better overall.

E.1 R2D2 experiments

Network architecture. R2D2, and our Retrace variants, use the 3-layer convolutional network from DQN [Mnih et al., 2015], followed by an LSTM with 512 hidden units, which then feeds into a dueling architecture of size 512 [Wang et al., 2016]. Like the original R2D2, the LSTM receives the reward and one-hot action vector from the previous time step as inputs.

Hyperparameters. The hyperparameters used for the R2D2 agents follow those of Kapturowski et al. [2019], and are reproduced in Table 1 for completeness.

Number of actors	256
Actor parameter update interval	400 environment steps
Sequence length m	80 (+ prefix of $l = 40$ for burn-in)
Replay buffer size	4×10^6 observations (10^5 part-overlapping sequences)
Priority exponent	0.9
Importance sampling exponent	0.6
Discount γ	0.997
Minibatch size	64
Optimiser	Adam [Kingma and Ba, 2015]
Optimiser settings	learning rate = 10^{-4} , $\epsilon = 10^{-3}$
Target network update interval	2500 updates
Value function rescaling	$h(x) = \text{sign}(x)(\sqrt{ x + 1} - 1) + \epsilon x$, $\epsilon = 10^{-3}$

Table 1: Hyperparameters values used in R2D2 experiments.

E.2 DQN experiments

Network architecture. The DQN, DoubleDQN, n -step and Retrace-based agents use the 3-layer convolutional network from DQN [Mnih et al., 2015], but unlike the R2D2 agents do not use an LSTM or dueling architecture.

Notice that the Retrace and C-trace agents are effectively using DoubleDQN-style updates due to the target probabilities not coming from the target network.

Hyperparameters. For sequential DQN-agents (n -step and Retrace) we performed a preliminary hyperparameter sweep to determine appropriate learning rates for n -step and Retrace updates. We swept over learning rates (0.00025, 0.0001, 0.00005, 0.00001) for both algorithms, and for n -step we jointly swept over two values for n (3 and 5). These were run on four Atari 2600 games (Alien, Amidar, Assault, Asterix), with the best performing hyperparameters for each method used for the Atari-57 experiments.

Interestingly, we found a small learning rate of 0.00001 worked best for both algorithms and that a larger $n = 5$ performed best for n -step.

Both algorithms used a maximum sequence length of 16. Due to shortness of the sequence length we use truncated trajectory corrections as described in the main text. Note that the truncation $\max(\Gamma, \gamma^N)$ is applied to each element of the sequence independently, therefore the value of N will begin at $N = 16$ for the first element and reduce to $N = 1$ for the final transition in the replay sequence.

F Further large-scale results

F.1 Detailed R2D2 results

We give further experimental results to complement the summary presented in the main paper; per-game training curves are given in Figure 9.

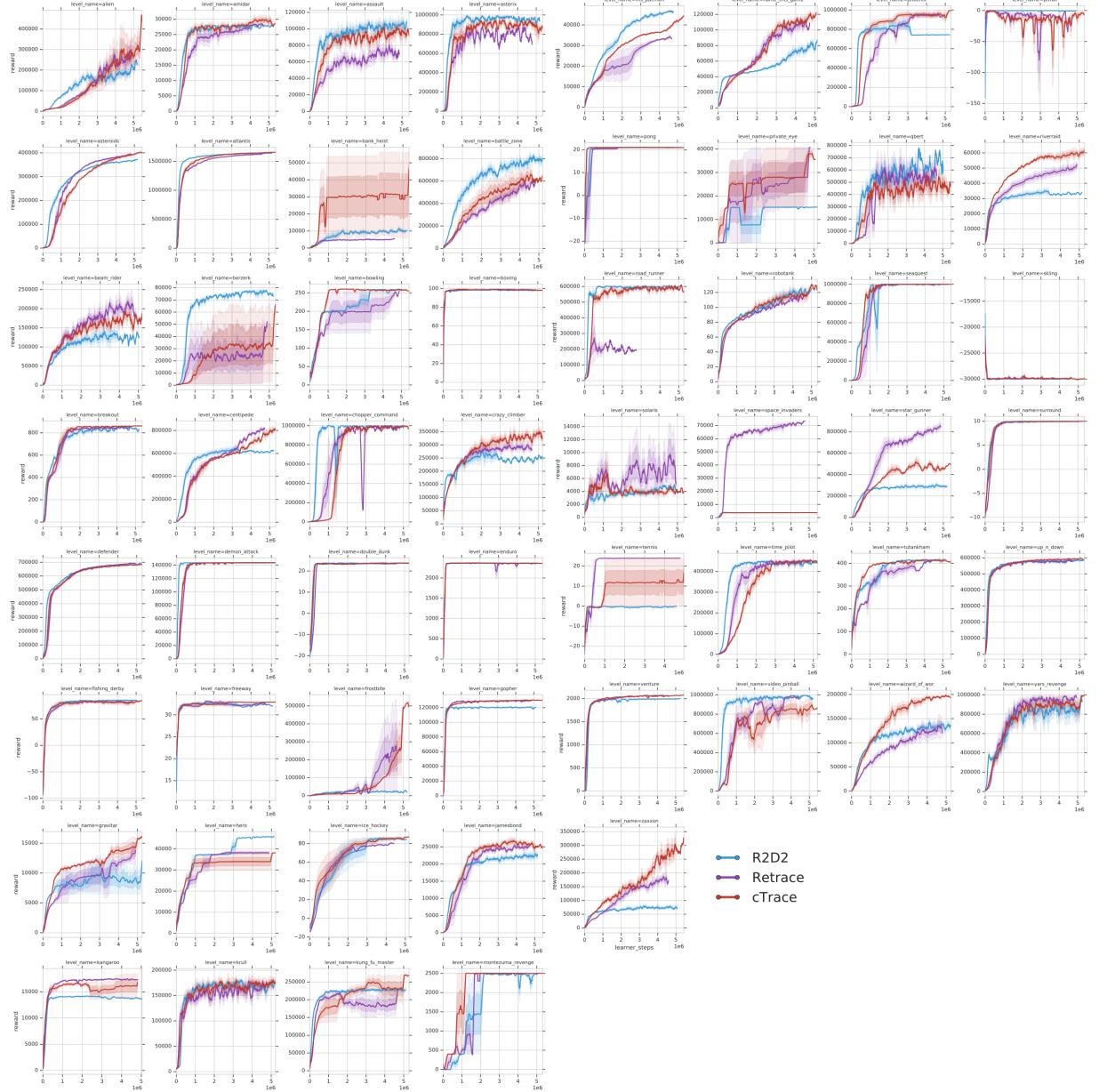


Figure 9: Training curves for 57 Atari games for R2D2 with n -step uncorrected returns (light blue), Retrace-R2D2 (black), and C-trace-R2D2 (red).

F.2 Detailed DQN results

We give further experimental results to complement the summary presented in the main paper. Results for varying the contraction hyperparameter are given in Figure 10, and per-game training curves for the main paper results are given in Figure 11.

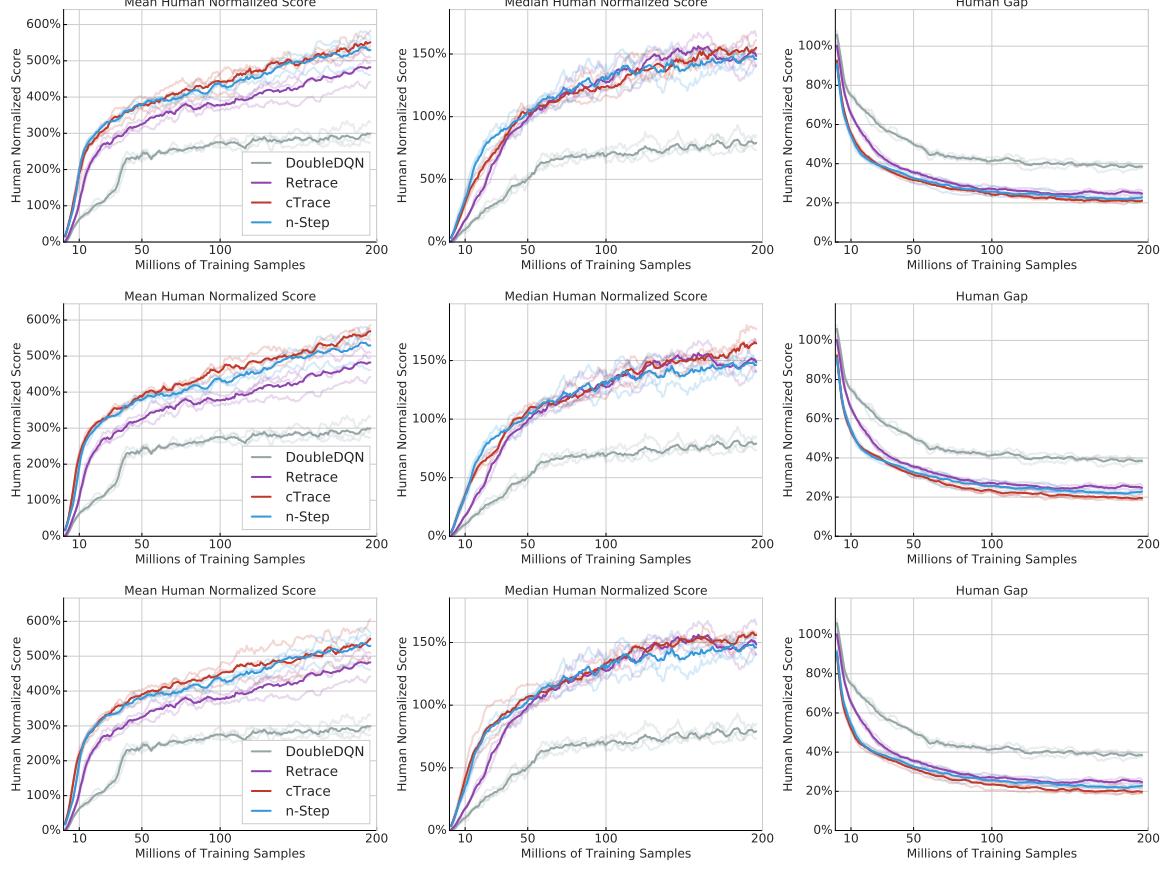


Figure 10: Atari-57 results for single-actor agent, as presented in the main text, but varying the C-trace contraction parameter: **(top)** γ^5 , **(center)** γ^7 , and **(bottom)** γ^{10} . Notice that due to its adaptation of α , C-trace is highly robust to the choice of contraction target.

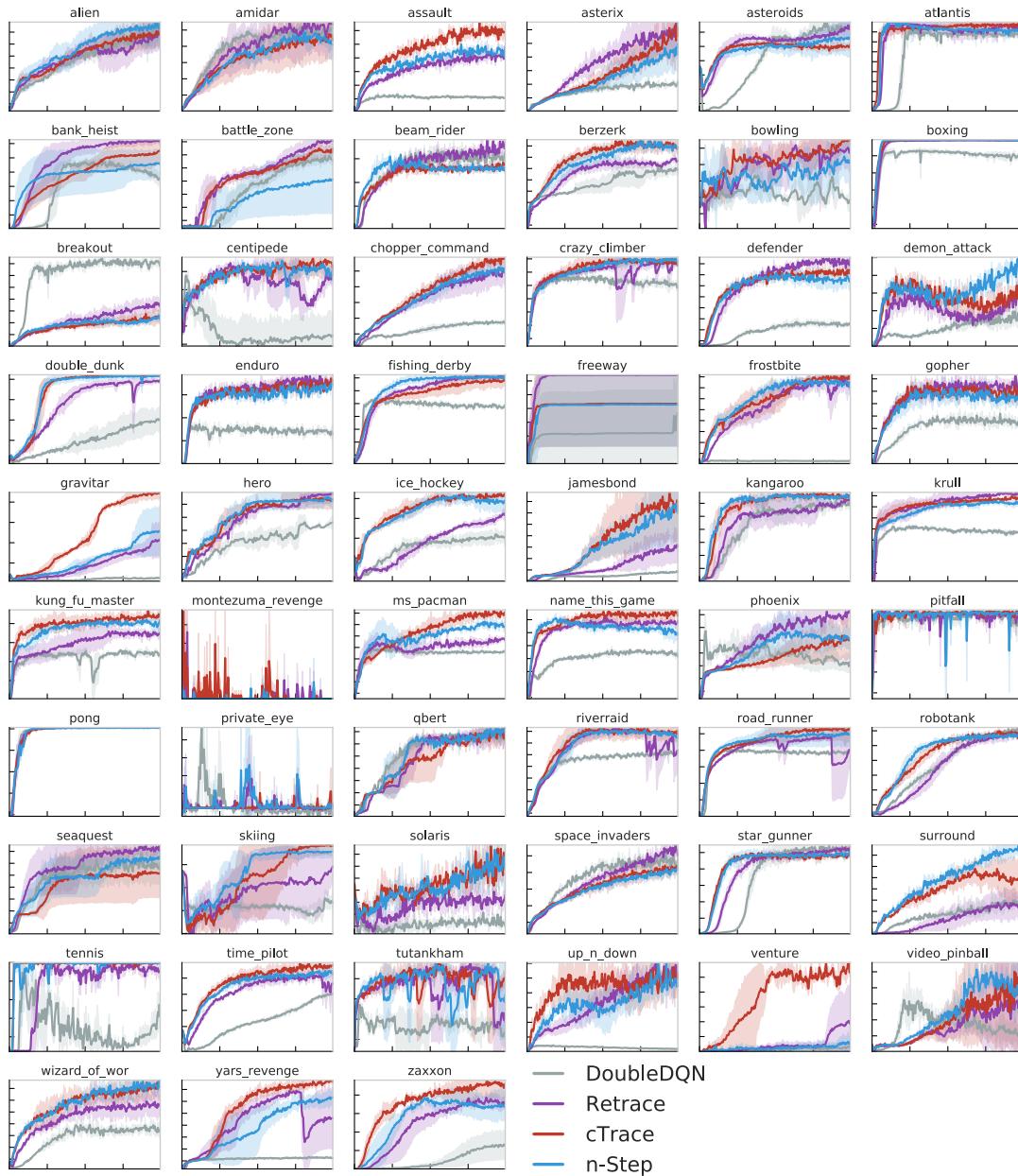


Figure 11: Training curves for 57 Atari games for Double DQN (grey), Double DQN with n -step uncorrected returns (light blue), Retrace-DQN (black), and C-trace-DQN (red).