

---

# Safe and efficient off-policy reinforcement learning

---

**Rémi Munos**  
munos@google.com  
Google DeepMind

**Thomas Stepleton**  
stepleton@google.com  
Google DeepMind

**Anna Harutyunyan**  
anna.harutyunyan@vub.ac.be  
Vrije Universiteit Brussel

**Marc G. Bellemare**  
bellemare@google.com  
Google DeepMind

## Abstract

In this work, we take a fresh look at some old and new algorithms for off-policy, return-based reinforcement learning. Expressing these in a common form, we derive a novel algorithm,  $\text{Retrace}(\lambda)$ , with three desired properties: (1) it has *low variance*; (2) it *safely* uses samples collected from any behaviour policy, whatever its degree of “off-policyness”; and (3) it is *efficient* as it makes the best use of samples collected from near on-policy behaviour policies. We analyze the contractive nature of the related operator under both off-policy policy evaluation and control settings and derive online sample-based algorithms. We believe this is the first return-based off-policy control algorithm converging a.s. to  $Q^*$  without the GLIE assumption (Greedy in the Limit with Infinite Exploration). As a corollary, we prove the convergence of Watkins’  $Q(\lambda)$ , which was an open problem since 1989. We illustrate the benefits of  $\text{Retrace}(\lambda)$  on a standard suite of Atari 2600 games.

One fundamental trade-off in reinforcement learning lies in the definition of the update target: should one estimate Monte Carlo returns or bootstrap from an existing Q-function? Return-based methods (where *return* refers to the sum of discounted rewards  $\sum_t \gamma^t r_t$ ) offer some advantages over value bootstrap methods: they are better behaved when combined with function approximation, and quickly propagate the fruits of exploration (Sutton, 1996). On the other hand, value bootstrap methods are more readily applied to off-policy data, a common use case. In this paper we show that *learning from returns need not be at cross-purposes with off-policy learning*.

We start from the recent work of Harutyunyan et al. (2016), who show that naive off-policy policy evaluation, without correcting for the “off-policyness” of a trajectory, still converges to the desired  $Q^\pi$  value function provided the behavior  $\mu$  and target  $\pi$  policies are not too far apart (the maximum allowed distance depends on the  $\lambda$  parameter). Their  $Q^\pi(\lambda)$  algorithm learns from trajectories generated by  $\mu$  simply by summing discounted off-policy corrected rewards at each time step. Unfortunately, the assumption that  $\mu$  and  $\pi$  are close is restrictive, as well as difficult to uphold in the control case, where the target policy is greedy with respect to the current Q-function. In that sense this algorithm is not *safe*: it does not handle the case of arbitrary “off-policyness”.

Alternatively, the Tree-backup (TB( $\lambda$ )) algorithm (Precup et al., 2000) tolerates arbitrary target/behavior discrepancies by scaling information (here called *traces*) from future temporal differences by the product of target policy probabilities. TB( $\lambda$ ) is not *efficient* in the “near on-policy” case (similar  $\mu$  and  $\pi$ ), though, as traces may be cut prematurely, blocking learning from full returns.

In this work, we express several off-policy, return-based algorithms in a common form. From this we derive an improved algorithm,  $\text{Retrace}(\lambda)$ , which is both *safe* and *efficient*, enjoying convergence guarantees for off-policy policy evaluation and – more importantly – for the control setting.

Retrace( $\lambda$ ) can learn from full returns retrieved from past policy data, as in the context of experience replay (Lin, 1993), which has returned to favour with advances in deep reinforcement learning (Mnih et al., 2015; Schaul et al., 2016). Off-policy learning is also desirable for exploration, since it allows the agent to deviate from the target policy currently under evaluation.

To the best of our knowledge, this is the first online return-based off-policy control algorithm which does not require the GLIE (Greedy in the Limit with Infinite Exploration) assumption (Singh et al., 2000). In addition, we provide as a corollary the first proof of convergence of Watkins’  $Q(\lambda)$  (see, e.g., Watkins, 1989; Sutton and Barto, 1998).

Finally, we illustrate the significance of Retrace( $\lambda$ ) in a deep learning setting by applying it to the suite of Atari 2600 games provided by the Arcade Learning Environment (Bellemare et al., 2013).

## 1 Notation

We consider an agent interacting with a Markov Decision Process  $(\mathcal{X}, \mathcal{A}, \gamma, P, r)$ .  $\mathcal{X}$  is a finite state space,  $\mathcal{A}$  the action space,  $\gamma \in [0, 1)$  the discount factor,  $P$  the transition function mapping state-action pairs  $(x, a) \in \mathcal{X} \times \mathcal{A}$  to distributions over  $\mathcal{X}$ , and  $r : \mathcal{X} \times \mathcal{A} \rightarrow [-R_{\max}, R_{\max}]$  is the reward function. For notational simplicity we will consider a finite action space, but the case of infinite – possibly continuous – action space can be handled by the Retrace( $\lambda$ ) algorithm as well. A policy  $\pi$  is a mapping from  $\mathcal{X}$  to a distribution over  $\mathcal{A}$ . A Q-function  $Q$  maps each state-action pair  $(x, a)$  to a value in  $\mathbb{R}$ ; in particular, the reward  $r$  is a Q-function. For a policy  $\pi$  we define the operator  $P^\pi$ :

$$(P^\pi Q)(x, a) := \sum_{x' \in \mathcal{X}} \sum_{a' \in \mathcal{A}} P(x' | x, a) \pi(a' | x') Q(x', a').$$

The value function for a policy  $\pi$ ,  $Q^\pi$ , describes the expected discounted sum of rewards associated with following  $\pi$  from a given state-action pair. Using operator notation, we write this as

$$Q^\pi := \sum_{t \geq 0} \gamma^t (P^\pi)^t r. \quad (1)$$

The *Bellman operator*  $\mathcal{T}^\pi$  for a policy  $\pi$  is defined as  $\mathcal{T}^\pi Q := r + \gamma P^\pi Q$  and its fixed point is  $Q^\pi$ , i.e.  $\mathcal{T}^\pi Q^\pi = Q^\pi = (I - \gamma P^\pi)^{-1} r$ . The *Bellman optimality operator* introduces a maximization over the set of policies:

$$\mathcal{T}Q := r + \gamma \max_{\pi} P^\pi Q. \quad (2)$$

Its fixed point is  $Q^*$ , the unique *optimal value function* (Puterman, 1994). It is this quantity that we will seek to obtain when we talk about the “control setting”.

**Return-based Operators:** The  $\lambda$ -return extension (Sutton, 1988) of the Bellman operators considers exponentially weighted sums of  $n$ -steps returns:

$$\mathcal{T}_\lambda^\pi Q := (1 - \lambda) \sum_{n \geq 0} \lambda^n [(\mathcal{T}^\pi)^{n+1} Q] = Q + (I - \lambda \gamma P^\pi)^{-1} (\mathcal{T}^\pi Q - Q),$$

where  $\mathcal{T}^\pi Q - Q$  is the *Bellman residual* of  $Q$  for policy  $\pi$ . Examination of the above shows that  $Q^\pi$  is also the fixed point of  $\mathcal{T}_\lambda^\pi$ . At one extreme ( $\lambda = 0$ ) we have the Bellman operator  $\mathcal{T}_{\lambda=0}^\pi Q = \mathcal{T}^\pi Q$ , while at the other ( $\lambda = 1$ ) we have the policy evaluation operator  $\mathcal{T}_{\lambda=1}^\pi Q = Q^\pi$  which can be estimated using Monte Carlo methods (Sutton and Barto, 1998). Intermediate values of  $\lambda$  trade off estimation bias with sample variance (Kearns and Singh, 2000).

We seek to evaluate a *target policy*  $\pi$  using trajectories drawn from a *behaviour policy*  $\mu$ . If  $\pi = \mu$ , we are *on-policy*; otherwise, we are *off-policy*. We will consider trajectories of the form:

$$x_0 = x, a_0 = a, r_0, x_1, a_1, r_1, x_2, a_2, r_2, \dots$$

with  $a_t \sim \mu(\cdot | x_t)$ ,  $r_t = r(x_t, a_t)$  and  $x_{t+1} \sim P(\cdot | x_t, a_t)$ . We denote by  $\mathcal{F}_t$  this sequence up to time  $t$ , and write  $\mathbb{E}_\mu$  the expectation with respect to both  $\mu$  and the MDP transition probabilities. Throughout, we write  $\|\cdot\|$  for supremum norm.

## 2 Off-Policy Algorithms

We are interested in two related off-policy learning problems. In the *policy evaluation* setting, we are given a fixed policy  $\pi$  whose value  $Q^\pi$  we wish to estimate from sample trajectories drawn from a behaviour policy  $\mu$ . In the *control* setting, we consider a sequence of policies that depend on our own sequence of Q-functions (such as  $\varepsilon$ -greedy policies), and seek to approximate  $Q^*$ .

The general operator that we consider for comparing several return-based off-policy algorithms is:

$$\mathcal{R}Q(x, a) := Q(x, a) + \mathbb{E}_\mu \left[ \sum_{t \geq 0} \gamma^t \left( \prod_{s=1}^t c_s \right) (r_t + \gamma \mathbb{E}_\pi Q(x_{t+1}, \cdot) - Q(x_t, a_t)) \right], \quad (3)$$

for some non-negative coefficients  $(c_s)$ , where we write  $\mathbb{E}_\pi Q(x, \cdot) := \sum_a \pi(a|x)Q(x, a)$  and define  $(\prod_{s=1}^t c_s) = 1$  when  $t = 0$ . By extension of the idea of eligibility traces (Sutton and Barto, 1998), we informally call the coefficients  $(c_s)$  the *traces* of the operator.

**Importance sampling (IS):**  $c_s = \frac{\pi(a_s|x_s)}{\mu(a_s|x_s)}$ . Importance sampling is the simplest way to correct for the discrepancy between  $\mu$  and  $\pi$  when learning from off-policy returns (Precup et al., 2000, 2001; Geist and Scherrer, 2014). The off-policy correction uses the product of the likelihood ratios between  $\pi$  and  $\mu$ . Notice that  $\mathcal{R}Q$  defined in (3) with this choice of  $(c_s)$  yields  $Q^\pi$  for any  $Q$ . For  $Q = 0$  we recover the basic IS estimate  $\sum_{t \geq 0} \gamma^t (\prod_{s=1}^t c_s) r_t$ , thus (3) can be seen as a variance reduction technique (with a baseline  $Q$ ). It is well known that IS estimates can suffer from large – even possibly infinite – variance (mainly due to the variance of the product  $\frac{\pi(a_1|x_1)}{\mu(a_1|x_1)} \dots \frac{\pi(a_t|x_t)}{\mu(a_t|x_t)}$ ), which has motivated further variance reduction techniques such as in (Mahmood and Sutton, 2015; Mahmood et al., 2015; Hallak et al., 2015).

**Off-policy  $Q^\pi(\lambda)$  and  $Q^*(\lambda)$ :**  $c_s = \lambda$ . A recent alternative proposed by Harutyunyan et al. (2016) introduces an off-policy correction based on a  $Q$ -baseline (instead of correcting the probability of the sample path like in IS). This approach, called  $Q^\pi(\lambda)$  and  $Q^*(\lambda)$  for policy evaluation and control, respectively, corresponds to the choice  $c_s = \lambda$ . It offers the advantage of avoiding the blow-up of the variance of the product of ratios encountered with IS. Interestingly, this operator contracts around  $Q^\pi$  provided that  $\mu$  and  $\pi$  are sufficiently close to each other. Defining  $\varepsilon := \max_x \|\pi(\cdot|x) - \mu(\cdot|x)\|_1$  the level of “off-policyness”, the authors prove that the operator defined by (3) with  $c_s = \lambda$  is a contraction mapping around  $Q^\pi$  for  $\lambda < \frac{1-\gamma}{\gamma\varepsilon}$ , and around  $Q^*$  for the worst case of  $\lambda < \frac{1-\gamma}{2\gamma}$ . Unfortunately,  $Q^\pi(\lambda)$  requires knowledge of  $\varepsilon$ , and the condition for  $Q^*(\lambda)$  is very conservative. Neither  $Q^\pi(\lambda)$ , nor  $Q^*(\lambda)$  are safe as they do not guarantee convergence for arbitrary  $\pi$  and  $\mu$ .

**Tree-backup, TB( $\lambda$ ):**  $c_s = \lambda \pi(a_s|x_s)$ . The TB( $\lambda$ ) algorithm of Precup et al. (2000) corrects for the target/behaviour discrepancy by multiplying each term of the sum by the product of target policy probabilities. The corresponding operator defines a contraction mapping for any policies  $\pi$  and  $\mu$ , which makes it a safe algorithm. However, this algorithm is not efficient in the near on-policy case (where  $\mu$  and  $\pi$  are similar) as it unnecessarily cuts the traces, preventing it to make use of full returns: indeed we need not discount stochastic on-policy transitions (as shown by Harutyunyan et al.’s results about  $Q^\pi$ ).

**Retrace( $\lambda$ ):**  $c_s = \lambda \min \left( 1, \frac{\pi(a_s|x_s)}{\mu(a_s|x_s)} \right)$ . Our contribution is an algorithm – Retrace( $\lambda$ ) – that takes the best of the three previous algorithms. Retrace( $\lambda$ ) uses an importance sampling ratio truncated at 1. Compared to IS, it does not suffer from the variance explosion of the product of IS ratios. Now, similarly to  $Q^\pi(\lambda)$  and unlike TB( $\lambda$ ), it does not cut the traces in the on-policy case, making it possible to benefit from the full returns. In the off-policy case, the traces are safely cut, similarly to TB( $\lambda$ ). In particular,  $\min \left( 1, \frac{\pi(a_s|x_s)}{\mu(a_s|x_s)} \right) \geq \pi(a_s|x_s)$ : Retrace( $\lambda$ ) does not cut the traces as much as TB( $\lambda$ ). In the subsequent sections, we will show the following:

- For any traces  $0 \leq c_s \leq \pi(a_s|x_s)/\mu(a_s|x_s)$  (thus including the Retrace( $\lambda$ ) operator), the return-based operator (3) is a  $\gamma$ -contraction around  $Q^\pi$ , for *arbitrary* policies  $\mu$  and  $\pi$
- In the control case (where  $\pi$  is replaced by a sequence of increasingly greedy policies) the online Retrace( $\lambda$ ) algorithm converges a.s. to  $Q^*$ , without requiring the GLIE assumption.
- As a corollary, Watkins’s  $Q(\lambda)$  converges a.s. to  $Q^*$ .

	Definition of $c_s$	Estimation variance	Guaranteed convergence <sup>†</sup>	Use full returns (near on-policy)
Importance sampling	$\frac{\pi(a_s x_s)}{\mu(a_s x_s)}$	High	for any $\pi, \mu$	yes
$Q^\pi(\lambda)$	$\lambda$	Low	for $\pi$ close to $\mu$	yes
TB( $\lambda$ )	$\lambda\pi(a_s x_s)$	Low	for any $\pi, \mu$	no
Retrace( $\lambda$ )	$\lambda \min\left(1, \frac{\pi(a_s x_s)}{\mu(a_s x_s)}\right)$	Low	for any $\pi, \mu$	yes

Table 1: Properties of several algorithms defined in terms of the general operator given in (3).  
<sup>†</sup>Guaranteed convergence of the expected operator  $\mathcal{R}$ .

### 3 Analysis of Retrace( $\lambda$ )

We will in turn analyze both off-policy policy evaluation and control settings. We will show that  $\mathcal{R}$  is a contraction mapping in both settings (under a mild additional assumption for the control case).

#### 3.1 Policy Evaluation

Consider a fixed target policy  $\pi$ . For ease of exposition we consider a fixed behaviour policy  $\mu$ , noting that our result extends to the setting of sequences of behaviour policies  $(\mu_k : k \in \mathbb{N})$ .

Our first result states the  $\gamma$ -contraction of the operator (3) defined by any set of non-negative coefficients  $c_s = c_s(a_s, \mathcal{F}_s)$  (in order to emphasize that  $c_s$  can be a function of the whole history  $\mathcal{F}_s$ ) under the assumption that  $0 \leq c_s \leq \frac{\pi(a_s|x_s)}{\mu(a_s|x_s)}$ .

**Theorem 1.** *The operator  $\mathcal{R}$  defined by (3) has a unique fixed point  $Q^\pi$ . Furthermore, if for each  $a_s \in \mathcal{A}$  and each history  $\mathcal{F}_s$  we have  $c_s = c_s(a_s, \mathcal{F}_s) \in [0, \frac{\pi(a_s|x_s)}{\mu(a_s|x_s)}]$ , then for any  $Q$ -function  $Q$*

$$\|\mathcal{R}Q - Q^\pi\| \leq \gamma \|Q - Q^\pi\|.$$

The following lemma will be useful in proving Theorem 1 (proof in the appendix).

**Lemma 1.** *The difference between  $\mathcal{R}Q$  and its fixed point  $Q^\pi$  is*

$$\mathcal{R}Q(x, a) - Q^\pi(x, a) = \mathbb{E}_\mu \left[ \sum_{t \geq 1} \gamma^t \left( \prod_{i=1}^{t-1} c_i \right) \left( [\mathbb{E}_\pi[(Q - Q^\pi)(x_t, \cdot)] - c_t(Q - Q^\pi)(x_t, a_t)] \right) \right].$$

*Proof (Theorem 1).* The fact that  $Q^\pi$  is the fixed point of the operator  $\mathcal{R}$  is obvious from (3) since  $\mathbb{E}_{x_{t+1} \sim P(\cdot|x_t, a_t)} [r_t + \gamma \mathbb{E}_\pi Q^\pi(x_{t+1}, \cdot) - Q^\pi(x_t, a_t)] = (\mathcal{T}^\pi Q^\pi - Q^\pi)(x_t, a_t) = 0$ , since  $Q^\pi$  is the fixed point of  $\mathcal{T}^\pi$ . Now, from Lemma 1, and defining  $\Delta Q := Q - Q^\pi$ , we have

$$\begin{aligned} \mathcal{R}Q(x, a) - Q^\pi(x, a) &= \sum_{t \geq 1} \gamma^t \mathbb{E}_{\substack{x_{1:t} \\ a_{1:t}}} \left[ \left( \prod_{i=1}^{t-1} c_i \right) \left( [\mathbb{E}_\pi \Delta Q(x_t, \cdot) - c_t \Delta Q(x_t, a_t)] \right) \right] \\ &= \sum_{t \geq 1} \gamma^t \mathbb{E}_{\substack{x_{1:t} \\ a_{1:t-1}}} \left[ \left( \prod_{i=1}^{t-1} c_i \right) \left( [\mathbb{E}_\pi \Delta Q(x_t, \cdot) - \mathbb{E}_{a_t} [c_t(a_t, \mathcal{F}_t) \Delta Q(x_t, a_t) | \mathcal{F}_t]] \right) \right] \\ &= \sum_{t \geq 1} \gamma^t \mathbb{E}_{\substack{x_{1:t} \\ a_{1:t-1}}} \left[ \left( \prod_{i=1}^{t-1} c_i \right) \sum_b (\pi(b|x_t) - \mu(b|x_t) c_t(b, \mathcal{F}_t)) \Delta Q(x_t, b) \right]. \end{aligned}$$

Now since  $\pi(a|x_t) - \mu(a|x_t) c_t(b, \mathcal{F}_t) \geq 0$ , we have that  $\mathcal{R}Q(x, a) - Q^\pi(x, a) = \sum_{y,b} w_{y,b} \Delta Q(y, b)$ , i.e. a linear combination of  $\Delta Q(y, b)$  weighted by non-negative coefficients:

$$w_{y,b} := \sum_{t \geq 1} \gamma^t \mathbb{E}_{\substack{x_{1:t} \\ a_{1:t-1}}} \left[ \left( \prod_{i=1}^{t-1} c_i \right) (\pi(b|x_t) - \mu(b|x_t) c_t(b, \mathcal{F}_t)) \mathbb{I}\{x_t = y\} \right].$$

The sum of those coefficients is:

$$\begin{aligned}
\sum_{y,b} w_{y,b} &= \sum_{t \geq 1} \gamma^t \mathbb{E}_{x_{1:t}} \left[ \left( \prod_{i=1}^{t-1} c_i \right) \sum_b (\pi(b|x_t) - \mu(b|x_t) c_t(b, \mathcal{F}_t)) \right] \\
&= \sum_{t \geq 1} \gamma^t \mathbb{E}_{x_{1:t}} \left[ \left( \prod_{i=1}^{t-1} c_i \right) \mathbb{E}_{a_t} [1 - c_t(a_t, \mathcal{F}_t) | \mathcal{F}_t] \right] = \sum_{t \geq 1} \gamma^t \mathbb{E}_{x_{1:t}} \left[ \left( \prod_{i=1}^{t-1} c_i \right) (1 - c_t) \right] \\
&= \mathbb{E}_\mu \left[ \sum_{t \geq 1} \gamma^t \left( \prod_{i=1}^{t-1} c_i \right) - \sum_{t \geq 1} \gamma^t \left( \prod_{i=1}^t c_i \right) \right] = \gamma C - (C - 1),
\end{aligned}$$

where  $C := \mathbb{E}_\mu \left[ \sum_{t \geq 0} \gamma^t \left( \prod_{i=1}^t c_i \right) \right]$ . Since  $C \geq 1$ , we have that  $\sum_{y,b} w_{y,b} \leq \gamma$ . Thus  $\mathcal{R}Q(x, a) - Q^\pi(x, a)$  is a sub-convex combination of  $\Delta Q(y, b)$  weighted by non-negative coefficients  $w_{y,b}$  which sum to (at most)  $\gamma$ , thus  $\mathcal{R}$  is a  $\gamma$ -contraction mapping around  $Q^\pi$ .  $\square$

**Remark 1.** Notice that the coefficient  $C$  in the proof of Theorem 1 depends on  $(x, a)$ . If we write  $\eta(x, a) := 1 - (1 - \gamma) \mathbb{E}_\mu \left[ \sum_{t \geq 0} \gamma^t \left( \prod_{s=1}^t c_s \right) \right]$ , then we have shown that

$$|\mathcal{R}Q(x, a) - Q^\pi(x, a)| \leq \eta(x, a) \|Q - Q^\pi\|.$$

Thus  $\eta(x, a) \in [0, \gamma]$  is a  $(x, a)$ -specific contraction coefficient, which is  $\gamma$  when  $c_1 = 0$  (the trace is cut immediately) and can be close to zero when learning from full returns ( $\mathbb{E}_\mu[c_t] \approx 1$  for all  $t$ ).

### 3.2 Control

In the control setting, the single target policy  $\pi$  is replaced by a sequence of policies  $(\pi_k)$  which depend on  $(Q_k)$ . While most prior work has focused on strictly greedy policies, here we consider the larger class of *increasingly greedy* sequences. We now make this notion precise.

**Definition 1.** We say that a sequence of policies  $(\pi_k : k \in \mathbb{N})$  is *increasingly greedy w.r.t. a sequence*  $(Q_k : k \in \mathbb{N})$  of  $Q$ -functions if the following property holds for all  $k$ :  $P^{\pi_{k+1}} Q_{k+1} \geq P^{\pi_k} Q_{k+1}$ .

Intuitively, this means that each  $\pi_{k+1}$  is at least as greedy as the previous policy  $\pi_k$  for  $Q_{k+1}$ . Many natural sequences of policies are increasingly greedy, including  $\varepsilon_k$ -greedy policies (with non-increasing  $\varepsilon_k$ ) and softmax policies (with non-increasing temperature). See proofs in the appendix.

We will assume that  $c_s = c_s(a_s, \mathcal{F}_s) = c(a_s, x_s)$  is Markovian, in the sense that it depends on  $x_s, a_s$  (as well as the policies  $\pi$  and  $\mu$ ) only but not on the full past history. This allows us to define the (sub)-probability transition operator

$$(P^{c\mu} Q)(x, a) := \sum_{x'} \sum_{a'} p(x'|x, a) \mu(a'|x') c(a', x') Q(x', a').$$

Finally, an additional requirement to the convergence in the control case, we assume that  $Q_0$  satisfies  $\mathcal{T}^{\pi_0} Q_0 \geq Q_0$  (this can be achieved by a pessimistic initialization  $Q_0 = -R_{MAX}/(1 - \gamma)$ ).

**Theorem 2.** Consider an arbitrary sequence of behaviour policies  $(\mu_k)$  (which may depend on  $(Q_k)$ ) and a sequence of target policies  $(\pi_k)$  that are increasingly greedy w.r.t. the sequence  $(Q_k)$ :

$$Q_{k+1} = \mathcal{R}_k Q_k,$$

where the return operator  $\mathcal{R}_k$  is defined by (3) for  $\pi_k$  and  $\mu_k$  and a Markovian  $c_s = c(a_s, x_s) \in [0, \frac{\pi_k(a_s|x_s)}{\mu_k(a_s|x_s)}]$ . Assume the target policies  $\pi_k$  are  $\varepsilon_k$ -away from the greedy policies w.r.t.  $Q_k$ , in the sense that  $\mathcal{T}^{\pi_k} Q_k \geq \mathcal{T} Q_k - \varepsilon_k \|Q_k\| e$ , where  $e$  is the vector with 1-components. Further suppose that  $\mathcal{T}^{\pi_0} Q_0 \geq Q_0$ . Then for any  $k \geq 0$ ,

$$\|Q_{k+1} - Q^*\| \leq \gamma \|Q_k - Q^*\| + \varepsilon_k \|Q_k\|.$$

In consequence, if  $\varepsilon_k \rightarrow 0$  then  $Q_k \rightarrow Q^*$ .

*Sketch of Proof (The full proof is in the appendix).* Using  $P^{c\mu_k}$ , the Retrace( $\lambda$ ) operator rewrites

$$\mathcal{R}_k Q = Q + \sum_{t \geq 0} \gamma^t (P^{c\mu_k})^t (\mathcal{T}^{\pi_k} Q - Q) = Q + (I - \gamma P^{c\mu_k})^{-1} (\mathcal{T}^{\pi_k} Q - Q).$$

We now lower- and upper-bound the term  $Q_{k+1} - Q^*$ .

**Upper bound on  $Q_{k+1} - Q^*$ .** We prove that  $Q_{k+1} - Q^* \leq A_k(Q_k - Q^*)$  with  $A_k := \gamma(I - \gamma P^{c\mu_k})^{-1} [P^{\pi_k} - P^{c\mu_k}]$ . Since  $c_t \in [0, \frac{\pi(a_t|x_t)}{\mu(a_t|x_t)}]$  we deduce that  $A_k$  has non-negative elements, whose sum over each row, is at most  $\gamma$ . Thus

$$Q_{k+1} - Q^* \leq \gamma \|Q_k - Q^*\| e. \quad (4)$$

**Lower bound on  $Q_{k+1} - Q^*$ .** Using the fact that  $\mathcal{T}^{\pi_k} Q_k \geq \mathcal{T}^{\pi^*} Q_k - \varepsilon_k \|Q_k\| e$  we have

$$\begin{aligned} Q_{k+1} - Q^* &\geq Q_{k+1} - \mathcal{T}^{\pi_k} Q_k + \gamma P^{\pi^*} (Q_k - Q^*) - \gamma \varepsilon_k \|Q_k\| e \\ &= \gamma P^{c\mu_k} (I - \gamma P^{c\mu_k})^{-1} (\mathcal{T}^{\pi_k} Q_k - Q_k) + \gamma P^{\pi^*} (Q_k - Q^*) - \varepsilon_k \|Q_k\| e. \end{aligned} \quad (5)$$

**Lower bound on  $\mathcal{T}^{\pi_k} Q_k - Q_k$ .** Since the sequence  $(\pi_k)$  is increasingly greedy w.r.t.  $(Q_k)$ , we have

$$\begin{aligned} \mathcal{T}^{\pi_{k+1}} Q_{k+1} - Q_{k+1} &\geq \mathcal{T}^{\pi_k} Q_{k+1} - Q_{k+1} = r + (\gamma P^{\pi_k} - I) \mathcal{R}_k Q_k \\ &= B_k (\mathcal{T}^{\pi_k} Q_k - Q_k), \end{aligned} \quad (6)$$

where  $B_k := \gamma [P^{\pi_k} - P^{c\mu_k}] (I - \gamma P^{c\mu_k})^{-1}$ . Since  $P^{\pi_k} - P^{c\mu_k}$  and  $(I - \gamma P^{c\mu_k})^{-1}$  are non-negative matrices, so is  $B_k$ . Thus  $\mathcal{T}^{\pi_k} Q_k - Q_k \geq B_{k-1} B_{k-2} \dots B_0 (\mathcal{T}^{\pi_0} Q_0 - Q_0) \geq 0$ , since we assumed  $\mathcal{T}^{\pi_0} Q_0 - Q_0 \geq 0$ . Thus, (5) implies that

$$Q_{k+1} - Q^* \geq \gamma P^{\pi^*} (Q_k - Q^*) - \varepsilon_k \|Q_k\| e.$$

Combining the above with (4) we deduce  $\|Q_{k+1} - Q^*\| \leq \gamma \|Q_k - Q^*\| + \varepsilon_k \|Q_k\|$ . When  $\varepsilon_k \rightarrow 0$ , we further deduce that  $Q_k$  are bounded, thus  $Q_k \rightarrow Q^*$ .  $\square$

### 3.3 Online algorithms

So far we have analyzed the contraction properties of the expected  $\mathcal{R}$  operators. We now describe online algorithms which can learn from sample trajectories. We analyze the algorithms in the *every visit* form (Sutton and Barto, 1998), which is the more practical generalization of the first-visit form. In this section, we will only consider the Retrace( $\lambda$ ) algorithm defined with the coefficient  $c = \lambda \min(1, \pi/\mu)$ . For that  $c$ , let us rewrite the operator  $P^{c\mu}$  as  $\lambda P^{\pi \wedge \mu}$ , where  $P^{\pi \wedge \mu} Q(x, a) := \sum_y \sum_b \min(\pi(b|y), \mu(b|y)) Q(y, b)$ , and write the Retrace operator  $\mathcal{R}Q = Q + (I - \lambda \gamma P^{\pi \wedge \mu})^{-1} (\mathcal{T}^{\pi} Q - Q)$ . We focus on the control case, noting that a similar (and simpler) result can be derived for policy evaluation.

**Theorem 3.** *Consider a sequence of sample trajectories, with the  $k^{th}$  trajectory  $x_0, a_0, r_0, x_1, a_1, r_1, \dots$  generated by following  $\mu_k$ :  $a_t \sim \mu_k(\cdot|x_t)$ . For each  $(x, a)$  along this trajectory, with  $s$  being the time of first occurrence of  $(x, a)$ , update*

$$Q_{k+1}(x, a) \leftarrow Q_k(x, a) + \alpha_k \sum_{t \geq s} \delta_t^{\pi_k} \sum_{j=s}^t \gamma^{t-j} \left( \prod_{i=j+1}^t c_i \right) \mathbb{I}\{x_j, a_j = x, a\}, \quad (7)$$

where  $\delta_t^{\pi_k} := r_t + \gamma \mathbb{E}_{\pi_k} Q_k(x_{t+1}, \cdot) - Q_k(x_t, a_t)$ ,  $\alpha_k = \alpha_k(x_s, a_s)$ . We consider the Retrace( $\lambda$ ) algorithm where  $c_i = \lambda \min(1, \frac{\pi(a_i|x_i)}{\mu(a_i|x_i)})$ . Assume that  $(\pi_k)$  are increasingly greedy w.r.t.  $(Q_k)$  and are each  $\varepsilon_k$ -away from the greedy policies  $(\pi_{Q_k})$ , i.e.  $\max_x \|\pi_k(\cdot|x) - \pi_{Q_k}(\cdot|x)\|_1 \leq \varepsilon_k$ , with  $\varepsilon_k \rightarrow 0$ . Assume that  $P^{\pi_k}$  and  $P^{\pi_k \wedge \mu_k}$  asymptotically commute:  $\lim_k \|P^{\pi_k} P^{\pi_k \wedge \mu_k} - P^{\pi_k \wedge \mu_k} P^{\pi_k}\| = 0$ . Assume further that (1) all states and actions are visited infinitely often:  $\sum_{t \geq 0} \mathbb{P}\{x_t, a_t = x, a\} \geq D > 0$ , (2) the sample trajectories are finite in terms of the second moment of their lengths  $T_k$ :  $\mathbb{E}_{\mu_k} T_k^2 < \infty$ , (3) the stepsizes obey the usual Robbins-Munro conditions. Then  $Q_k \rightarrow Q^*$  a.s.

The proof extends similar convergence proofs of TD( $\lambda$ ) by Bertsekas and Tsitsiklis (1996) and of optimistic policy iteration by Tsitsiklis (2003), and is provided in the appendix. Notice that compared to Theorem 2 we do not assume that  $\mathcal{T}^{\pi_0} Q_0 - Q_0 \geq 0$  here. However, we make the additional (rather technical) assumption that  $P^{\pi_k}$  and  $P^{\pi_k \wedge \mu_k}$  commute at the limit. This is satisfied for example when the probability assigned by the behavior policy  $\mu_k(\cdot|x)$  to the greedy action  $\pi_{Q_k}(x)$  is independent of  $x$ . Examples include  $\varepsilon$ -greedy policies, or more generally mixtures between the greedy policy  $\pi_{Q_k}$  and an arbitrary distribution  $\mu$  (see Lemma 5 in the appendix for the proof):

$$\mu_k(a|x) = \varepsilon \frac{\mu(a|x)}{1 - \mu(\pi_{Q_k}(x)|x)} \mathbb{I}\{a \neq \pi_{Q_k}(x)\} + (1 - \varepsilon) \mathbb{I}\{a = \pi_{Q_k}(x)\}. \quad (8)$$

Notice that the mixture coefficient  $\varepsilon$  needs not go to 0.

## 4 Discussion of the results

### 4.1 Choice of the trace coefficients $c_s$

Theorems 1 and 2 ensure convergence to  $Q^\pi$  and  $Q^*$  for any trace coefficient  $c_s \in [0, \frac{\pi(a_s|x_s)}{\mu(a_s|x_s)}]$ . However, to make the best choice of  $c_s$ , we need to consider the *speed* of convergence, which depends on both (1) the variance of the online estimate, which indicates how many online updates are required in a single iteration of  $\mathcal{R}$ , and (2) the contraction coefficient of  $\mathcal{R}$ .

**Variance:** The variance of the estimate strongly depends on the variance of the product trace  $(c_1 \dots c_t)$ , which is not an easy quantity to control in general, as the  $(c_s)$  are usually not independent. However, assuming independence and stationarity of  $(c_s)$ , we have that  $\mathbb{V}(\sum_t \gamma^t c_1 \dots c_t)$  is at least  $\sum_t \gamma^{2t} \mathbb{V}(c)^t$ , which is finite only if  $\mathbb{V}(c) < 1/\gamma^2$ . Thus, an important requirement for a numerically stable algorithm is for  $\mathbb{V}(c)$  to be as small as possible, and certainly no more than  $1/\gamma^2$ . This rules out importance sampling (for which  $c = \frac{\pi(a|x)}{\mu(a|x)}$ , and  $\mathbb{V}(c|x) = \sum_a \mu(a|x) (\frac{\pi(a|x)}{\mu(a|x)} - 1)^2$ , which may be larger than  $1/\gamma^2$  for some  $\pi$  and  $\mu$ ), and is the reason we choose  $c \leq 1$ .

**Contraction speed:** The contraction coefficient  $\eta \in [0, \gamma]$  of  $\mathcal{R}$  (see Remark 1) depends on how much the traces have been cut, and should be as small as possible (since it takes  $\log(1/\varepsilon)/\log(1/\eta)$  iterations of  $\mathcal{R}$  to obtain an  $\varepsilon$ -approximation). It is smallest when the traces are not cut at all (i.e. if  $c_s = 1$  for all  $s$ ,  $\mathcal{R}$  is the policy evaluation operator which produces  $Q^\pi$  in a single iteration). Indeed, when the traces are cut, we do not benefit from learning from full returns (in the extreme,  $c_1 = 0$  and  $\mathcal{R}$  reduces to the (one step) Bellman operator with  $\eta = \gamma$ ).

A reasonable trade-off between low variance (when  $c_s$  are small) and high contraction speed (when  $c_s$  are large) is given by  $\text{Retrace}(\lambda)$ , for which we provide the convergence of the online algorithm.

If we relax the assumption that the trace is Markovian (in which case only the result for policy evaluation has been proven so far) we could trade off a low trace at some time for a possibly larger-than-1 trace at another time, as long as their product is less than 1. A possible choice could be

$$c_s = \lambda \min \left( \frac{1}{c_1 \dots c_{s-1}}, \frac{\pi(a_s|x_s)}{\mu(a_s|x_s)} \right). \quad (9)$$

### 4.2 Other topics of discussion

**No GLIE assumption.** The crucial point of Theorem 2 is that convergence to  $Q^*$  occurs for *arbitrary* behaviour policies. Thus the online result in Theorem 3 does not require the behaviour policies to become greedy in the limit with infinite exploration (i.e. GLIE assumption, Singh et al., 2000). We believe Theorem 3 provides the first convergence result to  $Q^*$  for a  $\lambda$ -return (with  $\lambda > 0$ ) algorithm that does not require this (hard to satisfy) assumption.

**Proof of Watkins'  $Q(\lambda)$ .** As a corollary of Theorem 3 when selecting our target policies  $\pi_k$  to be greedy w.r.t.  $Q_k$  (i.e.  $\varepsilon_k = 0$ ), we deduce that Watkins'  $Q(\lambda)$  (e.g., Watkins, 1989; Sutton and Barto, 1998) converges a.s. to  $Q^*$  (under the assumption that  $\mu_k$  commutes asymptotically with the greedy policies, which is satisfied for e.g.  $\mu_k$  defined by (8)). We believe this is the first such proof.

**Increasingly greedy policies** The assumption that the sequence of target policies  $(\pi_k)$  is increasingly greedy w.r.t. the sequence of  $(Q_k)$  is more general than just considering greedy policies w.r.t.  $(Q_k)$  (which is Watkins's  $Q(\lambda)$ ), and leads to more efficient algorithms. Indeed, using non-greedy target policies  $\pi_k$  may speed up convergence as the traces are not cut as frequently. Of course, in order to converge to  $Q^*$ , we eventually need the target policies (and not the behaviour policies, as mentioned above) to become greedy in the limit (i.e.  $\varepsilon_k \rightarrow 0$  as defined in Theorem 2).

**Comparison to  $Q^\pi(\lambda)$ .** Unlike  $\text{Retrace}(\lambda)$ ,  $Q^\pi(\lambda)$  does not need to know the behaviour policy  $\mu$ . However, it fails to converge when  $\mu$  is far from  $\pi$ .  $\text{Retrace}(\lambda)$  uses its knowledge of  $\mu$  (for the chosen actions) to cut the traces and safely handle arbitrary policies  $\pi$  and  $\mu$ .

**Comparison to  $\text{TB}(\lambda)$ .** Similarly to  $Q^\pi(\lambda)$ ,  $\text{TB}(\lambda)$  does not need the knowledge of the behaviour policy  $\mu$ . But as a consequence,  $\text{TB}(\lambda)$  is not able to benefit from possible near on-policy situations, cutting traces unnecessarily when  $\pi$  and  $\mu$  are close.

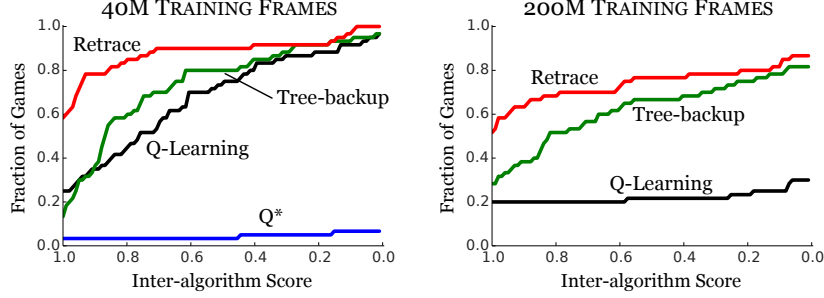


Figure 1: Inter-algorithm score distribution for  $\lambda$ -return ( $\lambda = 1$ ) variants and Q-Learning ( $\lambda = 0$ ).

**Estimating the behavior policy.** In the case  $\mu$  is unknown, it is reasonable to build an estimate  $\hat{\mu}$  from observed samples and use  $\hat{\mu}$  instead of  $\mu$  in the definition of the trace coefficients  $c_s$ . This may actually even lead to a better estimate, as analyzed by Li et al. (2015).

**Continuous action space.** Let us mention that Theorems 1 and 2 extend to the case of (measurable) continuous or infinite action spaces. The trace coefficients will make use of the densities  $\min(1, d\pi/d\mu)$  instead of the probabilities  $\min(1, \pi/\mu)$ . This is not possible with TB( $\lambda$ ).

**Open questions include:** (1) Removing the technical assumption that  $P^{\pi_k}$  and  $P^{\pi_k \wedge \mu_k}$  asymptotically commute, (2) Relaxing the Markov assumption in the control case in order to allow trace coefficients  $c_s$  of the form (9).

## 5 Experimental Results

To validate our theoretical results, we employ Retrace( $\lambda$ ) in an experience replay (Lin, 1993) setting, where sample transitions are stored within a large but bounded *replay memory* and subsequently replayed as if they were new experience. Naturally, older data in the memory is usually drawn from a policy which differs from the current policy, offering an excellent point of comparison for the algorithms presented in Section 2.

Our agent adapts the DQN architecture of Mnih et al. (2015) to replay short sequences from the memory (details in the appendix) instead of single transitions. The Q-function target update for a sample sequence  $x_t, a_t, r_t, \dots, x_{t+k}$  is

$$\Delta Q(x_t, a_t) = \sum_{s=t}^{t+k-1} \gamma^{s-t} \left( \prod_{i=t+1}^s c_i \right) [r(x_s, a_s) + \gamma \mathbb{E}_{\pi} Q(x_{s+1}, \cdot) - Q(x_s, a_s)].$$

We compare our algorithms' performance on 60 different Atari 2600 games in the Arcade Learning Environment (Bellemare et al., 2013) using Bellemare et al.'s inter-algorithm score distribution. Inter-algorithm scores are normalized so that 0 and 1 respectively correspond to the worst and best score for a particular game, within the set of algorithms under comparison. If  $g \in \{1, \dots, 60\}$  is a game and  $z_{g,a}$  the inter-algorithm score on  $g$  for algorithm  $a$ , then the score distribution function is  $f(x) := |\{g : z_{g,a} \geq x\}|/60$ . Roughly, a strictly higher curve corresponds to a better algorithm.

Across values of  $\lambda$ ,  $\lambda = 1$  performs best, save for  $Q^*(\lambda)$  where  $\lambda = 0.5$  obtains slightly superior performance. However, is highly sensitive to the choice of  $\lambda$  (see Figure 1, left, and Table 2 in the appendix). Both Retrace( $\lambda$ ) and TB( $\lambda$ ) achieve dramatically higher performance than Q-Learning early on and maintain their advantage throughout. Compared to TB( $\lambda$ ), Retrace( $\lambda$ ) offers a narrower but still marked advantage, being the best performer on 30 games; TB( $\lambda$ ) claims 15 of the remainder. Per-game details are given in the appendix.

**Conclusion.** Retrace( $\lambda$ ) can be seen as an algorithm that automatically adjusts – efficiently and safely – the length of the return to the degree of “off-policy” of any available data.

**Acknowledgments.** The authors thank Daan Wierstra, Nicolas Heess, Hado van Hasselt, Ziyu Wang, David Silver, Audrunas Gruslys, Georg Ostrovski, Hubert Soyer, and others at Google DeepMind for their very useful feedback on this work.

## References

- Bellemare, M. G., Naddaf, Y., Veness, J., and Bowling, M. (2013). The Arcade Learning Environment: An evaluation platform for general agents. *Journal of Artificial Intelligence Research*, 47:253–279.
- Bertsekas, D. P. and Tsitsiklis, J. N. (1996). *Neuro-Dynamic Programming*. Athena Scientific.
- Geist, M. and Scherrer, B. (2014). Off-policy learning with eligibility traces: A survey. *The Journal of Machine Learning Research*, 15(1):289–333.
- Hallak, A., Tamar, A., Munos, R., and Mannor, S. (2015). Generalized emphatic temporal difference learning: Bias-variance analysis. *arXiv:1509.05172*.
- Harutyunyan, A., Bellemare, M. G., Stepleton, T., and Munos, R. (2016).  $Q(\lambda)$  with off-policy corrections.
- Kearns, M. J. and Singh, S. P. (2000). Bias-variance error bounds for temporal difference updates. In *Conference on Computational Learning Theory*, pages 142–147.
- Li, L., Munos, R., and Szepesvari, C. (2015). Toward minimax off-policy value estimation. In *Proceedings of the 18th International Conference on Artificial Intelligence and Statistics (AISTATS)*.
- Lin, L. (1993). Scaling up reinforcement learning for robot control. In *Machine Learning: Proceedings of the Tenth International Conference*, pages 182–189.
- Mahmood, A. R. and Sutton, R. S. (2015). Off-policy learning based on weighted importance sampling with linear computational complexity. In *Conference on Uncertainty in Artificial Intelligence*.
- Mahmood, A. R., Yu, H., White, M., and Sutton, R. S. (2015). Emphatic temporal-difference learning. *arXiv:1507.01569*.
- Mnih, V., Badia, A. P., Mirza, M., Graves, A., Lillicrap, T. P., Harley, T., Silver, D., and Kavukcuoglu, K. (2016). Asynchronous methods for deep reinforcement learning. In *Proceedings of the International Conference on Machine Learning*.
- Mnih, V., Kavukcuoglu, K., Silver, D., Rusu, A. A., Veness, J., Bellemare, M. G., Graves, A., Riedmiller, M., Fidjeland, A. K., Ostrovski, G., et al. (2015). Human-level control through deep reinforcement learning. *Nature*, 518(7540):529–533.
- Precup, D., Sutton, R. S., and Dasgupta, S. (2001). Off-policy temporal-difference learning with function approximation. In *International Conference on Machine Learning*, pages 417–424.
- Precup, D., Sutton, R. S., and Singh, S. (2000). Eligibility traces for off-policy policy evaluation. In *Proceedings of the Seventeenth International Conference on Machine Learning*.
- Puterman, M. L. (1994). *Markov Decision Processes: Discrete Stochastic Dynamic Programming*. John Wiley & Sons, Inc., New York, NY, USA.
- Schaul, T., Quan, J., Antonoglou, I., and Silver, D. (2016). Prioritized experience replay. In *International Conference on Learning Representations*.
- Singh, S., Jaakkola, T., Littman, M. L., and Szepesvári, C. (2000). Convergence results for single-step on-policy reinforcement-learning algorithms. *Machine Learning*, 38(3):287–308.
- Sutton, R. and Barto, A. (1998). *Reinforcement learning: An introduction*. Cambridge Univ Press.
- Sutton, R. S. (1988). Learning to predict by the methods of temporal differences. *Machine learning*, 3(1):9–44.
- Sutton, R. S. (1996). Generalization in reinforcement learning: Successful examples using sparse coarse coding. In *Advances in Neural Information Processing Systems 8*.
- Tsitsiklis, J. N. (2003). On the convergence of optimistic policy iteration. *Journal of Machine Learning Research*, 3:59–72.
- Watkins, C. J. C. H. (1989). *Learning from Delayed Rewards*. PhD thesis, King’s College, Cambridge, UK.

## A Proof of Lemma 1

*Proof (Lemma 1).* Let  $\Delta Q := Q - Q^\pi$ . We begin by rewriting (3):

$$\mathcal{R}Q(x, a) = \sum_{t \geq 0} \gamma^t \mathbb{E}_\mu \left[ \left( \prod_{s=1}^t c_s \right) \left( r_t + \gamma \left[ \mathbb{E}_\pi Q(x_{t+1}, \cdot) - c_{t+1} Q(x_{t+1}, a_{t+1}) \right] \right) \right].$$

Since  $Q^\pi$  is the fixed point of  $\mathcal{R}$ , we have

$$Q^\pi(x, a) = \mathcal{R}Q^\pi(x, a) = \sum_{t \geq 0} \gamma^t \mathbb{E}_\mu \left[ \left( \prod_{s=1}^t c_s \right) \left( r_t + \gamma \left[ \mathbb{E}_\pi Q^\pi(x_{t+1}, \cdot) - c_{t+1} Q^\pi(x_{t+1}, a_{t+1}) \right] \right) \right],$$

from which we deduce that

$$\begin{aligned} \mathcal{R}Q(x, a) - Q^\pi(x, a) &= \sum_{t \geq 0} \gamma^t \mathbb{E}_\mu \left[ \left( \prod_{s=1}^t c_s \right) \left( \gamma \left[ \mathbb{E}_\pi \Delta Q(x_{t+1}, \cdot) - c_{t+1} \Delta Q(x_{t+1}, a_{t+1}) \right] \right) \right] \\ &= \sum_{t \geq 1} \gamma^t \mathbb{E}_\mu \left[ \left( \prod_{s=1}^{t-1} c_s \right) \left( \left[ \mathbb{E}_\pi \Delta Q(x_t, \cdot) - c_t \Delta Q(x_t, a_t) \right] \right) \right]. \end{aligned}$$

□

## B Increasingly greedy policies

Recall the definition of an increasingly greedy sequence of policies.

**Definition 2.** We say that a sequence of policies  $(\pi_k)$  is increasingly greedy w.r.t. a sequence of functions  $(Q_k)$  if the following property holds for all  $k$ :

$$P^{\pi_{k+1}} Q_{k+1} \geq P^{\pi_k} Q_{k+1}.$$

It is obvious to see that this property holds if all policies  $\pi_k$  are greedy w.r.t.  $Q_k$ . Indeed in such case,  $\mathcal{T}^{\pi_{k+1}} Q_{k+1} = \mathcal{T} Q_{k+1} \geq \mathcal{T}^{\pi_k} Q_{k+1}$  for any  $\pi$ .

We now prove that this property holds for  $\varepsilon_k$ -greedy policies (with non-increasing  $(\varepsilon_k)$ ) as well as soft-max policies (with non-decreasing  $(\beta_k)$ ), as stated in the two lemmas below.

Of course not all policies satisfy this property (a counter-example being  $\pi_k(a|x) := \arg \min_{a'} Q_k(x, a')$ ).

**Lemma 2.** Let  $(\varepsilon_k)$  be a non-increasing sequence. Then the sequence of policies  $(\pi_k)$  which are  $\varepsilon_k$ -greedy w.r.t. the sequence of functions  $(Q_k)$  is increasingly greedy w.r.t. that sequence.

*Proof.* From the definition of an  $\varepsilon$ -greedy policy we have:

$$\begin{aligned} P^{\pi_{k+1}} Q_{k+1}(x, a) &= \sum_y p(y|x, a) \left[ (1 - \varepsilon_{k+1}) \max_b Q_{k+1}(y, b) + \varepsilon_{k+1} \frac{1}{A} \sum_b Q_{k+1}(y, b) \right] \\ &\geq \sum_y p(y|x, a) \left[ (1 - \varepsilon_k) \max_b Q_{k+1}(y, b) + \varepsilon_k \frac{1}{A} \sum_b Q_{k+1}(y, b) \right] \\ &\geq \sum_y p(y|x, a) \left[ (1 - \varepsilon_k) Q_{k+1}(y, \arg \max_b Q_k(y, b)) + \varepsilon_k \frac{1}{A} \sum_b Q_{k+1}(y, b) \right] \\ &= P^{\pi_k} Q_{k+1}, \end{aligned}$$

where we used the fact that  $\varepsilon_{k+1} \leq \varepsilon_k$ . □

**Lemma 3.** Let  $(\beta_k)$  be a non-decreasing sequence of soft-max parameters. Then the sequence of policies  $(\pi_k)$  which are soft-max (with parameter  $\beta_k$ ) w.r.t. the sequence of functions  $(Q_k)$  is increasingly greedy w.r.t. that sequence.

*Proof.* For any  $Q$  and  $y$ , define  $\pi_\beta(b) = \frac{e^{\beta Q(y,b)}}{\sum_{b'} e^{\beta Q(y,b' )}}$  and  $f(\beta) = \sum_b \pi_\beta(b) Q(y, b)$ . Then we have

$$\begin{aligned} f'(\beta) &= \sum_b [\pi_\beta(b) Q(y, b) - \pi_\beta(b) \sum_{b'} \pi_\beta(b') Q(y, b')] Q(y, b) \\ &= \sum_b \pi_\beta(b) Q(y, b)^2 - \left( \sum_b \pi_\beta(b) Q(y, b) \right)^2 \\ &= \mathbb{V}_{b \sim \pi_\beta} [Q(y, b)] \geq 0. \end{aligned}$$

Thus  $\beta \mapsto f(\beta)$  is a non-decreasing function, and since  $\beta_{k+1} \geq \beta_k$ , we have

$$\begin{aligned} P^{\pi_{k+1}} Q_{k+1}(x, a) &= \sum_y p(y|x, a) \sum_b \frac{e^{\beta_{k+1} Q_{k+1}(y,b)}}{\sum_{b'} e^{\beta_{k+1} Q_{k+1}(y,b' )}} Q_{k+1}(y, b) \\ &\geq \sum_y p(y|x, a) \sum_b \frac{e^{\beta_k Q_{k+1}(y,b)}}{\sum_{b'} e^{\beta_k Q_{k+1}(y,b' )}} Q_{k+1}(y, b) \\ &= P^{\pi_k} Q_{k+1}(x, a). \end{aligned} \quad \square$$

## C Proof of Theorem 2

As mentioned in the main text, since  $c_s$  is Markovian, we can define the (sub)-probability transition operator

$$(P^{c\mu} Q)(x, a) := \sum_{x'} \sum_{a'} p(x'|x, a) \mu(a'|x') c(a', x') Q(x', a').$$

The Retrace( $\lambda$ ) operator then writes

$$\mathcal{R}_k Q = Q + \sum_{t \geq 0} \gamma^t (P^{c\mu_k})^t (\mathcal{T}^{\pi_k} Q - Q) = Q + (I - \gamma P^{c\mu_k})^{-1} (\mathcal{T}^{\pi_k} Q - Q).$$

*Proof.* We now lower- and upper-bound the term  $Q_{k+1} - Q^*$ .

**Upper bound on  $Q_{k+1} - Q^*$ .** Since  $Q_{k+1} = \mathcal{R}_k Q_k$ , we have

$$\begin{aligned} Q_{k+1} - Q^* &= Q_k - Q^* + (I - \gamma P^{c\mu_k})^{-1} [\mathcal{T}^{\pi_k} Q_k - Q_k] \\ &= (I - \gamma P^{c\mu_k})^{-1} [\mathcal{T}^{\pi_k} Q_k - Q_k + (I - \gamma P^{c\mu_k})(Q_k - Q^*)] \\ &= (I - \gamma P^{c\mu_k})^{-1} [\mathcal{T}^{\pi_k} Q_k - Q^* - \gamma P^{c\mu_k}(Q_k - Q^*)] \\ &= (I - \gamma P^{c\mu_k})^{-1} [\mathcal{T}^{\pi_k} Q_k - \mathcal{T} Q^* - \gamma P^{c\mu_k}(Q_k - Q^*)] \\ &\leq (I - \gamma P^{c\mu_k})^{-1} [\gamma P^{\pi_k}(Q_k - Q^*) - \gamma P^{c\mu_k}(Q_k - Q^*)] \\ &= \gamma (I - \gamma P^{c\mu_k})^{-1} [P^{\pi_k} - P^{c\mu_k}] (Q_k - Q^*), \\ &= A_k (Q_k - Q^*), \end{aligned} \tag{10}$$

where  $A_k := \gamma (I - \gamma P^{c\mu_k})^{-1} [P^{\pi_k} - P^{c\mu_k}]$ .

Now let us prove that  $A_k$  has non-negative elements, whose sum over each row is at most  $\gamma$ . Let  $e$  be the vector with 1-components. By rewriting  $A_k$  as  $\gamma \sum_{t \geq 0} \gamma^t (P^{c\mu_k})^t (P^{\pi_k} - P^{c\mu_k})$  and noticing that

$$(P^{\pi_k} - P^{c\mu_k})e(x, a) = \sum_{x'} \sum_{a'} p(x'|x, a) [\pi_k(a'|x') - c(a', x') \mu_k(a'|x')] \geq 0, \tag{11}$$

it is clear that all elements of  $A_k$  are non-negative. We have

$$\begin{aligned} A_k e &= \gamma \sum_{t \geq 0} \gamma^t (P^{c\mu_k})^t [P^{\pi_k} - P^{c\mu_k}] e \\ &= \gamma \sum_{t \geq 0} \gamma^t (P^{c\mu_k})^t e - \sum_{t \geq 0} \gamma^{t+1} (P^{c\mu_k})^{t+1} e \\ &= e - (1 - \gamma) \sum_{t \geq 0} \gamma^t (P^{c\mu_k})^t e \\ &\leq \gamma e, \end{aligned} \tag{12}$$

(since  $\sum_{t \geq 0} \gamma^t (P^{c\mu_k})^t e \geq e$ ). Thus  $A_k$  has non-negative elements, whose sum over each row, is at most  $\gamma$ . We deduce from (10) that  $Q_{k+1} - Q^*$  is upper-bounded by a sub-convex combination of components of  $Q_k - Q^*$ ; the sum of their coefficients is at most  $\gamma$ . Thus

$$Q_{k+1} - Q^* \leq \gamma \|Q_k - Q^*\| e. \quad (13)$$

**Lower bound on  $Q_{k+1} - Q^*$ .** We have

$$\begin{aligned} Q_{k+1} &= Q_k + (I - \gamma P^{c\mu_k})^{-1} (\mathcal{T}^{\pi_k} Q_k - Q_k) \\ &= Q_k + \sum_{i \geq 0} \gamma^i (P^{c\mu_k})^i (\mathcal{T}^{\pi_k} Q_k - Q_k) \\ &= \mathcal{T}^{\pi_k} Q_k + \sum_{i \geq 1} \gamma^i (P^{c\mu_k})^i (\mathcal{T}^{\pi_k} Q_k - Q_k) \\ &= \mathcal{T}^{\pi_k} Q_k + \gamma P^{c\mu_k} (I - \gamma P^{c\mu_k})^{-1} (\mathcal{T}^{\pi_k} Q_k - Q_k). \end{aligned} \quad (14)$$

Now, from the definition of  $\varepsilon_k$  we have  $\mathcal{T}^{\pi_k} Q_k \geq \mathcal{T} Q_k - \varepsilon_k \|Q_k\| \geq \mathcal{T}^{\pi^*} Q_k - \varepsilon_k \|Q_k\|$ , thus

$$\begin{aligned} Q_{k+1} - Q^* &= Q_{k+1} - \mathcal{T}^{\pi_k} Q_k + \mathcal{T}^{\pi_k} Q_k - \mathcal{T}^{\pi^*} Q_k + \mathcal{T}^{\pi^*} Q_k - \mathcal{T}^{\pi^*} Q^* \\ &\geq Q_{k+1} - \mathcal{T}^{\pi_k} Q_k + \gamma P^{\pi^*} (Q_k - Q^*) - \varepsilon_k \|Q_k\| e \end{aligned}$$

Using (14) we derive the lower bound:

$$Q_{k+1} - Q^* \geq \gamma P^{c\mu_k} (I - \gamma P^{c\mu_k})^{-1} (\mathcal{T}^{\pi_k} Q_k - Q_k) + \gamma P^{\pi^*} (Q_k - Q^*) - \varepsilon_k \|Q_k\|. \quad (15)$$

**Lower bound on  $\mathcal{T}^{\pi_k} Q_k - Q_k$ .** By hypothesis,  $(\pi_k)$  is increasingly greedy w.r.t.  $(Q_k)$ , thus

$$\begin{aligned} \mathcal{T}^{\pi_{k+1}} Q_{k+1} - Q_{k+1} &\geq \mathcal{T}^{\pi_k} Q_{k+1} - Q_{k+1} \\ &= \mathcal{T}^{\pi_k} \mathcal{R}_k Q_k - \mathcal{R}_k Q_k \\ &= r + (\gamma P^{\pi_k} - I) \mathcal{R}_k Q_k \\ &= r + (\gamma P^{\pi_k} - I) [Q_k + (I - \gamma P^{c\mu_k})^{-1} (\mathcal{T}^{\pi_k} Q_k - Q_k)] \\ &= \mathcal{T}^{\pi_k} Q_k - Q_k + (\gamma P^{\pi_k} - I) (I - \gamma P^{c\mu_k})^{-1} (\mathcal{T}^{\pi_k} Q_k - Q_k) \\ &= \gamma [P^{\pi_k} - P^{c\mu_k}] (I - \gamma P^{c\mu_k})^{-1} (\mathcal{T}^{\pi_k} Q_k - Q_k) \\ &= B_k (\mathcal{T}^{\pi_k} Q_k - Q_k), \end{aligned} \quad (16)$$

where  $B_k := \gamma [P^{\pi_k} - P^{c\mu_k}] (I - \gamma P^{c\mu_k})^{-1}$ . Since  $P^{\pi_k} - P^{c\mu_k}$  has non-negative elements (as proven in (11)) as well as  $(I - \gamma P^{c\mu_k})^{-1}$ , then  $B_k$  has non-negative elements as well. Thus

$$\mathcal{T}^{\pi_k} Q_k - Q_k \geq B_{k-1} B_{k-2} \dots B_0 (\mathcal{T}^{\pi_0} Q_0 - Q_0) \geq 0,$$

since we assumed  $\mathcal{T}^{\pi_0} Q_0 - Q_0 \geq 0$ . Thus (15) implies that

$$Q_{k+1} - Q^* \geq \gamma P^{\pi^*} (Q_k - Q^*) - \varepsilon_k \|Q_k\|.$$

and combining the above with (13) we deduce

$$\|Q_{k+1} - Q^*\| \leq \gamma \|Q_k - Q^*\| + \varepsilon_k \|Q_k\|.$$

Now assume that  $\varepsilon_k \rightarrow 0$ . We first deduce that  $Q_k$  is bounded. Indeed as soon as  $\varepsilon_k < (1 - \gamma)/2$ , we have

$$\|Q_{k+1}\| \leq \|Q^*\| + \gamma \|Q_k - Q^*\| + \frac{1 - \gamma}{2} \|Q_k\| \leq (1 + \gamma) \|Q^*\| + \frac{1 + \gamma}{2} \|Q_k\|.$$

Thus  $\limsup \|Q_k\| \leq \frac{1 + \gamma}{1 - (1 + \gamma)/2} \|Q^*\|$ . Since  $Q_k$  is bounded, we deduce that  $\limsup Q_k = Q^*$ .  $\square$

## D Proof of Theorem 3

We first prove convergence of the general online algorithm.

**Theorem 4.** *Consider the algorithm*

$$Q_{k+1}(x, a) = (1 - \alpha_k(x, a))Q_k(x, a) + \alpha_k(x, a)(\mathcal{R}_k Q_k(x, a) + \omega_k(x, a) + v_k(x, a)), \quad (17)$$

and assume that (1)  $\omega_k$  is a centered,  $\mathcal{F}_k$ -measurable noise term of bounded variance, and (2)  $v_k$  is bounded from above by  $\theta_k(\|Q_k\| + 1)$ , where  $(\theta_k)$  is a random sequence that converges to 0 a.s. Then, under the same assumptions as in Theorem 3, we have that  $Q_k \rightarrow Q^*$  almost surely.

*Proof.* We write  $\mathcal{R}$  for  $\mathcal{R}_k$ . Let us prove the result in three steps.

**Upper bound on  $\mathcal{R}Q_k - Q^*$ .** The first part of the proof is similar to the proof of (13), so we have

$$\mathcal{R}Q_k - Q^* \leq \gamma\|Q_k - Q^*\|e. \quad (18)$$

**Lower bound on  $\mathcal{R}Q_k - Q^*$ .** Again, similarly to (15) we have

$$\begin{aligned} \mathcal{R}Q_k - Q^* &\geq \gamma\lambda P^{\pi_k \wedge \mu_k} (I - \gamma\lambda P^{\pi_k \wedge \mu_k})^{-1} (\mathcal{T}^{\pi_k} Q_k - Q_k) \\ &\quad + \gamma P^{\pi^*} (Q_k - Q^*) - \varepsilon_k \|Q_k\|. \end{aligned} \quad (19)$$

**Lower-bound on  $\mathcal{T}^{\pi_k} Q_k - Q_k$ .** Since the sequence of policies  $(\pi_k)$  is increasingly greedy w.r.t.  $(Q_k)$ , we have

$$\begin{aligned} \mathcal{T}^{\pi_{k+1}} Q_{k+1} - Q_{k+1} &\geq \mathcal{T}^{\pi_k} Q_{k+1} - Q_{k+1} \\ &= (1 - \alpha_k) \mathcal{T}^{\pi_k} Q_k + \alpha_k \mathcal{T}^{\pi_k} (\mathcal{R}Q_k + \omega_k + v_k) - Q_{k+1} \\ &= (1 - \alpha_k) (\mathcal{T}^{\pi_k} Q_k - Q_k) + \alpha_k [\mathcal{T}^{\pi_k} \mathcal{R}Q_k - \mathcal{R}Q_k + \omega'_k + v'_k], \end{aligned} \quad (20)$$

where  $\omega'_k := (\gamma P^{\pi_k} - I)\omega_k$  and  $v'_k := (\gamma P^{\pi_k} - I)v_k$ . It is easy to see that both  $\omega'_k$  and  $v'_k$  continue to satisfy the assumptions on  $\omega_k$ , and  $v_k$ . Now, from the definition of the  $\mathcal{R}$  operator, we have

$$\begin{aligned} \mathcal{T}^{\pi_k} \mathcal{R}Q_k - \mathcal{R}Q_k &= r + (\gamma P^{\pi_k} - I) \mathcal{R}Q_k \\ &= r + (\gamma P^{\pi_k} - I) [Q_k + (I - \gamma\lambda P^{\pi_k \wedge \mu_k})^{-1} (\mathcal{T}^{\pi_k} Q_k - Q_k)] \\ &= \mathcal{T}^{\pi_k} Q_k - Q_k + (\gamma P^{\pi_k} - I) (I - \gamma\lambda P^{\pi_k \wedge \mu_k})^{-1} (\mathcal{T}^{\pi_k} Q_k - Q_k) \\ &= \gamma (P^{\pi_k} - \lambda P^{\pi_k \wedge \mu_k}) (I - \gamma\lambda P^{\pi_k \wedge \mu_k})^{-1} (\mathcal{T}^{\pi_k} Q_k - Q_k). \end{aligned}$$

Using this equality into (20) and writing  $\xi_k := \mathcal{T}^{\pi_k} Q_k - Q_k$ , we have

$$\xi_{k+1} \geq (1 - \alpha_k) \xi_k + \alpha_k [B_k \xi_k + \omega'_k + v'_k], \quad (21)$$

where  $B_k := \gamma(P^{\pi_k} - \lambda P^{\pi_k \wedge \mu_k})(I - \gamma\lambda P^{\pi_k \wedge \mu_k})^{-1}$ . The matrix  $B_k$  is non-negative but may not be a contraction mapping (the sum of its components per row may be larger than 1). Thus we cannot directly apply Proposition 4.5 of Bertsekas and Tsitsiklis (1996). However, as we have seen in the proof of Theorem 2, the matrix  $A_k := \gamma(I - \gamma\lambda P^{\pi_k \wedge \mu_k})^{-1}(P^{\pi_k} - \lambda P^{\pi_k \wedge \mu_k})$  is a  $\gamma$ -contraction mapping. So now we relate  $B_k$  to  $A_k$  using our assumption that  $P^{\pi_k}$  and  $P^{\pi_k \wedge \mu_k}$  commute asymptotically, i.e.  $\|P^{\pi_k} P^{\pi_k \wedge \mu_k} - P^{\pi_k \wedge \mu_k} P^{\pi_k}\| = \eta_k$  with  $\eta_k \rightarrow 0$ . For any (sub)-transition matrices  $U$  and  $V$ , we have

$$\begin{aligned} U(I - \lambda\gamma V)^{-1} &= \sum_{t \geq 0} (\lambda\gamma)^t U V^t \\ &= \sum_{t \geq 0} (\lambda\gamma)^t \left[ \sum_{s=0}^{t-1} V^s (UV - VU) V^{t-s-1} + V^t U \right] \\ &= (I - \lambda\gamma V)^{-1} U + \sum_{t \geq 0} (\lambda\gamma)^t \sum_{s=0}^{t-1} V^s (UV - VU) V^{t-s-1}. \end{aligned}$$

Replacing  $U$  by  $P^{\pi_k}$  and  $V$  by  $P^{\pi_k \wedge \mu_k}$ , we deduce

$$\|B_k - A_k\| \leq \gamma \sum_{t \geq 0} t (\lambda\gamma)^t \eta_k = \gamma \frac{1}{(1 - \lambda\gamma)^2} \eta_k.$$

Thus, from (21),

$$\xi_{k+1} \geq (1 - \alpha_k)\xi_k + \alpha_k [A_k \xi_k + \omega'_k + v''_k], \quad (22)$$

where  $v''_k := v'_k + \gamma \sum_{t \geq 0} t(\lambda\gamma)^t \eta_k \|\xi_k\|$  continues to satisfy the assumptions on  $v_k$  (since  $\eta_k \rightarrow 0$ ).

Now, let us define another sequence  $\xi'_k$  as follows:  $\xi'_0 = \xi_0$  and

$$\xi'_{k+1} = (1 - \alpha_k)\xi'_k + \alpha_k (A_k \xi'_k + \omega'_k + v''_k).$$

We can now apply Proposition 4.5 of Bertsekas and Tsitsiklis (1996) to the sequence  $(\xi'_k)$ . The matrices  $A_k$  are non-negative, and the sum of their coefficients per row is bounded by  $\gamma$ , see (12), thus  $A_k$  are  $\gamma$ -contraction mappings and have the same fixed point which is 0. The noise  $\omega'_k$  is centered and  $\mathcal{F}_k$ -measurable and satisfies the bounded variance assumption, and  $v''_k$  is bounded above by  $(1 + \gamma)\theta'_k(\|Q_k\| + 1)$  for some  $\theta'_k \rightarrow 0$ . Thus  $\lim_k \xi'_k = 0$  almost surely.

Now, it is straightforward to see that  $\xi_k \geq \xi'_k$  for all  $k \geq 0$ . Indeed by induction, let us assume that  $\xi_k \geq \xi'_k$ . Then

$$\begin{aligned} \xi_{k+1} &\geq (1 - \alpha_k)\xi_k + \alpha_k (A_k \xi_k + \omega'_k + v''_k) \\ &\geq (1 - \alpha_k)\xi'_k + \alpha_k (A_k \xi'_k + \omega'_k + v''_k) \\ &= \xi'_{k+1}, \end{aligned}$$

since all elements of the matrix  $A_k$  are non-negative. Thus we deduce that

$$\liminf_{k \rightarrow \infty} \xi_k \geq \lim_{k \rightarrow \infty} \xi'_k = 0 \quad (23)$$

**Conclusion.** Using (23) in (19) we deduce the lower bound:

$$\liminf_{k \rightarrow \infty} \mathcal{R}Q_k - Q^* \geq \liminf_{k \rightarrow \infty} \gamma P^{\pi^*}(Q_k - Q^*), \quad (24)$$

almost surely. Now combining with the upper bound (18) we deduce that

$$\|\mathcal{R}Q_k - Q^*\| \leq \gamma \|Q_k - Q^*\| + O(\varepsilon_k \|Q_k\|) + O(\xi_k).$$

The last two terms can be incorporated to the  $v_k(x, a)$  and  $\omega_k(x, a)$  terms, respectively; we thus again apply Proposition 4.5 of Bertsekas and Tsitsiklis (1996) to the sequence  $(Q_k)$  defined by (17) and deduce that  $Q_k \rightarrow Q^*$  almost surely.  $\square$

It remains to rewrite the update (7) in the form of (17), in order to apply Theorem 4.

Let  $z_{s,t}^k$  denote the accumulating trace (Sutton and Barto, 1998):

$$z_{s,t}^k := \sum_{j=s}^t \gamma^{t-j} \left( \prod_{i=j+1}^t c_i \right) \mathbb{I}\{(x_j, a_j) = (x_s, a_s)\}.$$

Let us write  $Q_{k+1}^o(x_s, a_s)$  to emphasize the online setting. Then (7) can be written as

$$Q_{k+1}^o(x_s, a_s) \leftarrow Q_k^o(x_s, a_s) + \alpha_k(x_s, a_s) \sum_{t \geq s} \delta_t^{\pi_k} z_{s,t}^k, \quad (25)$$

$$\delta_t^{\pi_k} := r_t + \gamma \mathbb{E}_{\pi_k} Q_k^o(x_{t+1}, \cdot) - Q_k^o(x_t, a_t),$$

Using our assumptions on finite trajectories, and  $c_i \leq 1$ , we can show that:

$$\mathbb{E} \left[ \sum_{t \geq s} z_{s,t}^k | \mathcal{F}_k \right] < \mathbb{E} [T_k^2 | \mathcal{F}_k] < \infty \quad (26)$$

where  $T_k$  denotes trajectory length. Now, let  $D_k := D_k(x_s, a_s) := \sum_{t \geq s} \mathbb{P}\{(x_t, a_t) = (x_s, a_s)\}$ . Then, using (26), we can show that the total update is bounded, and rewrite

$$\mathbb{E}_{\mu_k} \left[ \sum_{t \geq s} \delta_t^{\pi_k} z_{s,t}^k \right] = D_k(x_s, a_s) (\mathcal{R}Q_k(x_s, a_s) - Q(x_s, a_s)).$$

Finally, using the above, and writing  $\alpha_k = \alpha_k(x_s, a_s)$ , (25) can be rewritten in the desired form:

$$\begin{aligned} Q_{k+1}^o(x_s, a_s) &\leftarrow (1 - \tilde{\alpha}_k)Q_k^o(x_s, a_s) + \tilde{\alpha}_k(\mathcal{R}_k Q_k^o(x_s, a_s) + \omega_k(x_s, a_s) + v_k(x_s, a_s)), \quad (27) \\ \omega_k(x_s, a_s) &:= (D_k)^{-1} \left( \sum_{t \geq s} \delta_t^{\pi_k} z_{s,t}^k - \mathbb{E}_{\mu_k} \left[ \sum_{t \geq s} \delta_t^{\pi_k} z_{s,t}^k \right] \right), \\ v_k(x_s, a_s) &:= (\tilde{\alpha}_k)^{-1} (Q_{k+1}^o(x_s, a_s) - Q_{k+1}(x_s, a_s)), \\ \tilde{\alpha}_k &:= \alpha_k D_k. \end{aligned}$$

It can be shown that the variance of the noise term  $\omega_k$  is bounded, using (26) and the fact that the reward function is bounded. It follows from Assumptions 1-3 that the modified stepsize sequence  $(\tilde{\alpha}_k)$  satisfies the conditions of Assumption 1. The second noise term  $v_k(x_s, a_s)$  measures the difference between online iterates and the corresponding offline values, and can be shown to satisfy the required assumption analogously to the argument in the proof of Prop. 5.2 in Bertsekas and Tsitsiklis (1996). The proof relies on the eligibility coefficients (26) and rewards being bounded, the trajectories being finite, and the conditions on the stepsizes being satisfied.

We can thus apply Theorem 4 to (27), and conclude that the iterates  $Q_k^o \rightarrow Q^*$  as  $k \rightarrow \infty$ , w.p. 1.

## E Asymptotic commutativity of $P^{\pi_k}$ and $P^{\pi_k \wedge \mu_k}$

**Lemma 4.** *Let  $(\pi_k)$  and  $(\mu_k)$  two sequences of policies. If there exists  $\alpha$  such that for all  $x, a$ ,*

$$\min(\pi_k(a|x), \mu_k(a|x)) = \alpha \pi_k(a|x) + o(1), \quad (28)$$

*then the transition matrices  $P^{\pi_k}$  and  $P^{\pi_k \wedge \mu_k}$  asymptotically commute:  $\|P^{\pi_k} P^{\pi_k \wedge \mu_k} - P^{\pi_k \wedge \mu_k} P^{\pi_k}\| = o(1)$ .*

*Proof.* For any  $Q$ , we have

$$\begin{aligned} (P^{\pi_k} P^{\pi_k \wedge \mu_k})Q(x, a) &= \sum_y p(y|x, a) \sum_b \pi_k(b|y) \sum_z p(z|y, b) \sum_c (\pi_k \wedge \mu_k)(c|z) Q(z, c) \\ &= \alpha \sum_y p(y|x, a) \sum_b \pi_k(b|y) \sum_z p(z|y, b) \sum_c \pi_k(c|z) Q(z, c) + \|Q\|o(1) \\ &= \sum_y p(y|x, a) \sum_b (\pi_k \wedge \mu_k)(b|y) \sum_z p(z|y, b) \sum_c \pi_k(c|z) Q(z, c) + \|Q\|o(1) \\ &= (P^{\pi_k \wedge \mu_k} P^{\pi_k})Q(x, a) + \|Q\|o(1). \quad \square \end{aligned}$$

**Lemma 5.** *Let  $(\pi_{Q_k})$  a sequence of (deterministic) greedy policies w.r.t. a sequence  $(Q_k)$ . Let  $(\pi_k)$  a sequence of policies that are  $\varepsilon_k$  away from  $(\pi_{Q_k})$ , in the sense that, for all  $x$ ,*

$$\|\pi_k(\cdot|x) - \pi_{Q_k}(x)\|_1 := 1 - \pi_k(\pi_{Q_k}(x)|x) + \sum_{a \neq \pi_{Q_k}(x)} \pi_k(a|x) \leq \varepsilon_k.$$

*Let  $(\mu_k)$  a sequence of policies defined by:*

$$\mu_k(a|x) = \frac{\alpha \mu(a|x)}{1 - \mu(\pi_{Q_k}(x)|x)} \mathbb{I}\{a \neq \pi_{Q_k}(x)\} + (1 - \alpha) \mathbb{I}\{a = \pi_{Q_k}(x)\}, \quad (29)$$

*for some arbitrary policy  $\mu$  and  $\alpha \in [0, 1]$ . Assume  $\varepsilon_k \rightarrow 0$ . Then the transition matrices  $P^{\pi_k}$  and  $P^{\pi_k \wedge \mu_k}$  asymptotically commute.*

*Proof.* The intuition is that asymptotically  $\pi_k$  gets very close to the deterministic policy  $\pi_{Q_k}$ . In that case, the minimum distribution  $(\pi_k \wedge \mu_k)(\cdot|x)$  puts a mass close to  $1 - \alpha$  on the greedy action  $\pi_{Q_k}(x)$ , and no mass on other actions, thus  $(\pi_k \wedge \mu_k)$  gets very close to  $(1 - \alpha)\pi_k$ , and Lemma 4 applies (with multiplicative constant  $1 - \alpha$ ).

Indeed, from our assumption that  $\pi_k$  is  $\varepsilon$ -away from  $\pi_{Q_k}$  we have:

$$\pi_k(\pi_{Q_k}(x)|x) \geq 1 - \varepsilon_k, \text{ and } \pi_k(a \neq \pi_{Q_k}(x)|x) \leq \varepsilon_k.$$

We deduce that

$$\begin{aligned}
(\pi_k \wedge \mu_k)(\pi_{Q_k}(x)|x) &= \min(\pi_k(\pi_{Q_k}(x)|x), 1 - \alpha) \\
&= 1 - \alpha + O(\varepsilon_k) \\
&= (1 - \alpha)\pi_k(\pi_{Q_k}(x)|x) + O(\varepsilon_k),
\end{aligned}$$

and

$$\begin{aligned}
(\pi_k \wedge \mu_k)(a \neq \pi_{Q_k}(x)|x) &= O(\varepsilon_k) \\
&= (1 - \alpha)\pi_k(a|x) + O(\varepsilon_k).
\end{aligned}$$

Thus Lemma 4 applies (with a multiplicative constant  $1 - \alpha$ ) and  $P^{\pi_k}$  and  $P^{\pi_k \wedge \mu_k}$  asymptotically commute.  $\square$

## F Experimental Methods

Although our experiments’ learning problem closely matches the DQN setting used by Mnih et al. (2015) (i.e. single-thread off-policy learning with large replay memory), we conducted our trials in the multi-threaded, CPU-based framework of Mnih et al. (2016), obtaining ample result data from affordable CPU resources. Key differences from the DQN are as follows. Sixteen threads with private environment instances train simultaneously; each infers with and finds gradients w.r.t. a local copy of the network parameters; gradients then update a “master” parameter set and local copies are refreshed. Target network parameters are simply shared globally. Each thread has private replay memory holding 62,500 transitions (1/16<sup>th</sup> of DQN’s total replay capacity). The optimizer is unchanged from (Mnih et al., 2016): “Shared RMSprop” with step size annealing to 0 over  $3 \times 10^8$  environment frames (summed over threads). Exploration parameter ( $\varepsilon$ ) behaviour differs slightly: every 50,000 frames, threads switch randomly (probability 0.3, 0.4, and 0.3 respectively) between three schedules (anneal  $\varepsilon$  from 1 to 0.5, 0.1, or 0.01 over 250,000 frames), starting new schedules at the intermediate positions where they left old ones.<sup>1</sup>

Our experiments comprise 60 Atari 2600 games in ALE (Bellemare et al., 2013), with “life” loss treated as episode termination. The control, minibatched (64 transitions/minibatch) one-step Q-learning as in (Mnih et al., 2015), shows performance comparable to DQN in our multi-threaded setup. Retrace, TB, and  $Q^*$  runs use minibatches of four 16-step sequences (again 64 transitions/minibatch) and the current exploration policy as the target policy  $\pi$ . All trials clamp rewards into  $[-1, 1]$ . In the control, Q-function targets are clamped into  $[-1, 1]$  prior to gradient calculation; analogous quantities in the multi-step algorithms are clamped into  $[-1, 1]$ , then scaled (divided by) the sequence length. Coarse, then fine logarithmic parameter sweeps on the games *Asterix*, *Breakout*, *Enduro*, *Freeway*, *H.E.R.O.*, *Pong*, *Q\*bert*, and *Seaquest* yielded step sizes of 0.0000439 and 0.0000912, and RMSprop regularization parameters of 0.001 and 0.0000368, for control and multi-step algorithms respectively. Reported performance averages over four trials with different random seeds for each experimental configuration.

### F.1 Algorithmic Performance in Function of $\lambda$

We compared our algorithms for different values of  $\lambda$ , using the DQN score as a baseline. As before, for each  $\lambda$  we compute the inter-algorithm scores on a per-game basis. We then averaged the inter-algorithm scores across games to produce Table 2 (see also Figure 2 for a visual depiction). We first remark that Retrace always achieve a score higher than TB, demonstrating that it is efficient in the sense of Section 2. Next, we note that  $Q^*$  performs best for small values of  $\lambda$ , but begins to fail for values above  $\lambda = 0.5$ . In this sense, it is also not safe. This is particularly problematic as the safe threshold of  $\lambda$  is likely to be problem-dependent. Finally, there is no setting of  $\lambda$  for which Retrace performs particularly poorly; for high values of  $\lambda$ , it achieves close to the top score in most games. For Retrace( $\lambda$ ) it makes sense to use a values  $\lambda = 1$  (at least in deterministic environments) as the trace cutting effect required in off-policy learning is taken care of by the use of the  $\min(1, \pi/\mu)$  coefficient. On the contrary,  $Q^*(\lambda)$  only relies on a value  $\lambda < 1$  to take care of cutting traces for off-policy data.

<sup>1</sup>We evaluated a DQN-style single schedule for  $\varepsilon$ , but our multi-schedule method, similar to the one used by Mnih et al., yielded improved performance in our multi-threaded setting.

$\lambda$	DQN	TB	Retrace	$Q^*$
0.0	0.5071	0.5512	0.4288	0.4487
0.1	0.4752	0.2798	0.5046	0.651
0.3	0.3634	0.268	0.5159	0.7734
0.5	0.2409	0.4105	0.5098	0.8419
0.7	0.3712	0.4453	0.6762	0.5551
0.9	0.7256	0.7753	0.9034	0.02926
1.0	0.6839	0.8158	0.8698	0.04317

Table 2: Average inter-algorithm scores for each value of  $\lambda$ . The DQN scores are fixed across different  $\lambda$ , but the corresponding inter-algorithm scores varies depending on the worst and best performer within each  $\lambda$ .

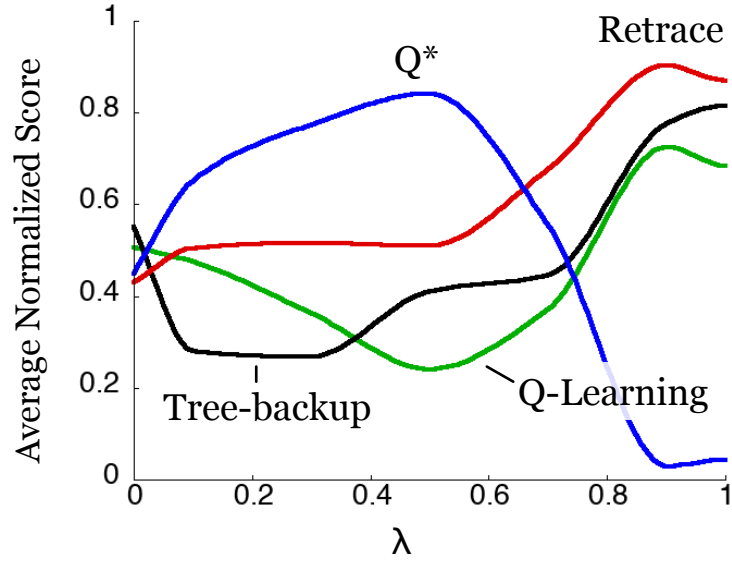


Figure 2: Average inter-algorithm scores for each value of  $\lambda$ . The DQN scores are fixed across different  $\lambda$ , but the corresponding inter-algorithm scores varies depending on the worst and best performer within each  $\lambda$ . **Note that average scores are not directly comparable across different values of  $\lambda$ .**

	Tree-backup( $\lambda$ )	Retrace( $\lambda$ )	DQN	Q*( $\lambda$ )
ALIEN	2508.62	<b>3109.21</b>	2088.81	154.35
AMIDAR	1221.00	<b>1247.84</b>	772.30	16.04
ASSAULT	7248.08	<b>8214.76</b>	1647.25	260.95
ASTERIX	<b>29294.76</b>	28116.39	10675.57	285.44
ASTEROIDS	1499.82	<b>1538.25</b>	1403.19	308.70
ATLANTIS	<b>2115949.75</b>	2110401.90	1712671.88	3667.18
BANK HEIST	<b>808.31</b>	797.36	549.35	1.70
BATTLE ZONE	22197.96	<b>23544.08</b>	21700.01	3278.93
BEAM RIDER	15931.60	<b>17281.24</b>	8053.26	621.40
BERZERK	967.29	<b>972.67</b>	627.53	247.80
BOWLING	40.96	<b>47.92</b>	37.82	15.16
BOXING	91.00	93.54	<b>95.17</b>	-29.25
BREAKOUT	288.71	298.75	<b>332.67</b>	1.21
CARNIVAL	<b>4691.73</b>	4633.77	4637.86	353.10
CENTIPEDE	1199.46	1715.95	1037.95	<b>3783.60</b>
CHOPPER COMMAND	6193.28	<b>6358.81</b>	5007.32	534.83
CRAZY CLIMBER	<b>115345.95</b>	114991.29	111918.64	1136.21
DEFENDER	32411.77	<b>33146.83</b>	13349.26	1838.76
DEMON ATTACK	68148.22	<b>79954.88</b>	8585.17	310.45
DOUBLE DUNK	<b>-1.32</b>	-6.78	-5.74	-23.63
ELEVATOR ACTION	1544.91	2396.05	<b>14607.10</b>	930.38
ENDURO	1115.00	<b>1216.47</b>	938.36	12.54
FISHING DERBY	22.22	<b>27.69</b>	15.14	-98.58
FREEWAY	<b>32.13</b>	32.13	31.07	9.86
FROSTBITE	960.30	935.42	<b>1124.60</b>	45.07
GOPHER	13666.33	<b>14110.94</b>	11542.46	50.59
GRAVITAR	30.18	29.04	<b>271.40</b>	13.14
H.E.R.O.	<b>25048.33</b>	21989.46	17626.90	12.48
ICE HOCKEY	<b>-3.84</b>	-5.08	-4.36	-15.68
JAMES BOND	560.88	641.51	<b>705.55</b>	21.71
KANGAROO	11755.01	<b>11896.25</b>	4101.92	178.23
KRULL	<b>9509.83</b>	9485.39	7728.66	429.26
KUNG-FU MASTER	25338.05	<b>26695.19</b>	17751.73	39.99
MONTEZUMA'S REVENGE	<b>0.79</b>	0.18	0.10	0.00
MS. PAC-MAN	2461.10	<b>3208.03</b>	2654.97	298.58
NAME THIS GAME	<b>11358.81</b>	11160.15	10098.85	1311.73
PHOENIX	13834.27	<b>15637.88</b>	9249.38	107.41
PITFALL	<b>-37.74</b>	-43.85	-392.63	-121.99
POOYAN	5283.69	<b>5661.92</b>	3301.69	98.65
PONG	<b>20.25</b>	20.20	19.31	-20.99
PRIVATE EYE	73.44	<b>87.36</b>	44.73	-147.49
Q*BERT	13617.24	<b>13700.25</b>	12412.85	114.84
RIVER RAID	14457.29	<b>15365.61</b>	10329.58	922.13
ROAD RUNNER	34396.52	32843.09	<b>50523.75</b>	418.62
ROBOTANK	36.07	41.18	<b>49.20</b>	5.77
SEAQUEST	3557.09	2914.00	<b>3869.30</b>	175.29
SKIING	-25055.94	-25235.75	-25254.43	<b>-24179.71</b>
SOLARIS	1178.05	1135.51	<b>1258.02</b>	674.58
SPACE INVADERS	<b>6096.21</b>	5623.34	2115.80	227.39
STAR GUNNER	66369.18	<b>74016.10</b>	42179.52	266.15
SURROUND	<b>-5.48</b>	-6.04	-8.17	-9.98
TENNIS	-1.73	-0.30	<b>13.67</b>	-7.37
TIME PILOT	8266.79	<b>8719.19</b>	8228.89	657.59
TUTANKHAM	164.54	<b>199.25</b>	167.22	2.68
UP AND DOWN	14976.51	<b>18747.40</b>	9404.95	530.59
VENTURE	10.75	22.84	<b>30.93</b>	0.09
VIDEO PINBALL	103486.09	<b>228283.79</b>	76691.75	6837.86
WIZARD OF WOR	7402.99	<b>8048.72</b>	612.52	189.43
YAR'S REVENGE	14581.65	<b>26860.57</b>	15484.03	1913.19
ZAXXON	12529.22	<b>15383.11</b>	8422.49	0.40
Times Best	16	30	12	2

Table 3: Final scores achieved by the different  $\lambda$ -return variants ( $\lambda = 1$ ). Highlights indicate high scores.