

1 We thank all reviewers and we will modify the paper to clarify each of the points raised, as discussed/clarified below.

2 **Similar points across reviewers:** (1) Re statistical significance of empirical results, our presentation was misleading
3 and instead we will provide confidence intervals (CIs) that clarify/quantify the statistical significance of our results.
4 I.e., the 95% CIs of all mean scores are relatively small for all operators; e.g., such CIs for cartpole are $< \pm 0.3$ for
5 each operator, much smaller than the differences in their mean scores. (2) Re results for constant β , we will expand
6 the discussion in Sec 4.5 (L338-341) noting constant β performs worse, and will provide the corresponding numerical
7 results; e.g., mean cartpole scores for $\beta = 1$ and $\beta \sim U[0, 1]$ are 187 compared with 191 for $\beta \sim U[0, 2]$. (3) Re defs
8 of terms, instead of providing a reference as in L92-L94, we will add such defs. E.g.: stochastic ordering (s.o.): r.v.
9 X is stochastically \leq to r.v. Y if $\mathbb{P}[X > z] \leq \mathbb{P}[Y > z], \forall z$; convex ordering (c.o.): r.v. X is $<$ r.v. Y under c.o. iff
10 $\mathbb{E}[f(X)] \leq \mathbb{E}[f(Y)], \forall$ convex functions f . (4) Re theoretical results establishing benefits of our approach, we will
11 provide some more intuition. For any (x, a) in eq (4) where the action a is not the optimal action, there will very often
12 (i.e. for many k) be instances where $V_k(x) > Q_k(x, a)$, so eq (4) will make $Q(x, a)$ (eventually) deviate more from
13 $V(x)$; OTOH, for action a s.t. $Q(x, a) = V(x)$, then $V_k(x) > Q_k(x, a)$ will only happen rarely, and thus eq (4) will
14 not affect the end value of $V(x)$. These both reflect the concepts of optimality preserving and action gap increasing.
15 Moreover, we observe that the multiplier in front of $V_k(x) - Q_k(x, a)$ (i.e. β_k) is desired to be large individually,
16 but its overall efforts should not be so large as to affect $V(x)$. We therefore introduce a family of RSOs, where β_k
17 is allowed to take on any value, but its average remains < 1 . Furthermore, we establish that greater variability in β_k
18 will lead to larger action gaps and that s.o. (c.o.) in β_k will lead to s.o. (c.o.) in the action gaps. (5) Re theoretical
19 proofs, because our RSOs introduce probabilistic elements on top of the original MDP, it is natural for us to employ
20 probabilistic arguments in our analysis/proofs. E.g., in L421-424, we use the lim sup and lim inf for set sequences for
21 the ease of derivations. The probabilistic nature of the problem also affords us the liberty to exploit weak convergence
22 limits (convergence in probability) to identify the limit of $V_k(x)$ after establishing a stronger a.s. convergence for $V_k(x)$.
23 More importantly, the stochastic nature of the problem allows us to consider s.o. (c.o.), which are common machinery
24 in probability theory and which we exploit to establish important orderings of performance among different RSOs.
25 While some of this may not be very familiar within the AI community, we believe these additions broaden the spectrum
26 of ideas and methodologies that can be exploited to help improve solutions of fundamental problems in RL and beyond.

27 **R1.** (1) Benefit of stochastic β_k is addressed by our theoretical results (Thms 3.2-3.4) and by our empirical results
28 demonstrating significant improvements under stochastic β_k . (2) L260 ff. were intended to note that distributions with
29 lower means and variances performed worse than $U[0, 2]$ in our experiments, which we will clarify/expand. Further
30 Thms 3.2-3.4 are intended to help find the best β_k sequence. (3) Indeed, the submission contained a typo in eq (9): π^{b_k}
31 in both inequalities should be a , which is the focus of the arguments that follow. (4) We should have explicitly stated
32 that $V_k(x) + f_k$ is uniformly upper bounded from the facts that the rewards are bounded functions and $\gamma \in (0, 1)$. (5)
33 In L421-L424, we establish the (right) inequality in eq. (9) by considering the limit superior and limit inferior of the
34 sequence of sets on which the probabilities are calculated. Therefore, we examine events that happen infinitely often for
35 the lim sup and all but finite exceptions for lim inf. (6) We view the problem of finding the maximally efficient operator
36 as one of finding a sequence of β_k that produces dominating performance. We believe that the statement is correct from
37 this viewpoint. But the statement does not discuss how the optimization should be conducted, where different methods
38 could have different implications. We will clarify these points. (7) The value function $V(x)$ indeed does not change,
39 but $Q(x, a)$ changes and this leads to a larger action gap. This should then lead to more efficient ways of ultimately
40 finding $V(x)$ via updating $Q(x, a)$, as indicated in refs [5] and [12].

41 **R2.** To address concerns of the role of randomness, we will expand L263-265 which notes the same trends were
42 observed when we varied exploration of the ϵ -greedy algorithm over a wide range of ϵ (even for deterministic operators).

43 **R3.** (1) We appreciate the comment on “robust” in other contexts, and will be more careful/clearer in our usage. (2) Re
44 proof of Thm 3.1, the inequality in L408 follows from $V_k(x) \geq Q_k(x, a), \forall a$, by defn of $V_k(x)$. The 3rd relation below
45 L411 follows directly from the defn of \mathcal{T}_B and the 4th relation below L411 is directly due to the order relation of \mathcal{T}_B
46 and \mathcal{T}_{β_k} . Lastly, regarding the a.s. convergence and its weak form, convergence in probability: By L417, we already
47 established that $V_k(x)$ converges a.s., so the purpose of the next paragraph is to identify the limit, as stated; For that
48 purpose, we only need to identify the limit of the weak convergence, since a.s. convergence naturally implies that $V_k(x)$
49 also weakly converges to the limit. (3) Re proof of Thm 3.2, it is quite standard to prove results by introducing a general
50 function with minimally restrictive properties (increasing) and then appropriately using this function and its properties;
51 the ordering $\mathbb{E}[f(\hat{Q}_{k+1})] \geq \mathbb{E}[f(\bar{Q}_{k+1})]$ leads to $\hat{Q}_k \geq_{st} \bar{Q}_k$ because of the properties of $f(\cdot)$ and the defn of s.o.
52 (\geq_{st}). Similarly, in the last part of the proof, $f(\cdot)$ is used to establish the desired s.o. of the action gap. (4) Re proof of
53 Thm 3.3, with the addition of the defn of c.o., the proof follows along similar lines to that of Thm 3.2 as stated. (5)
54 The proof of Thm 3.4 is not based on induction and is proven using a basic relation between variance and conditional
55 expectation, with the direct derivation establishing the general result for each k . (6) Q-learning provides convergence to
56 the values of all $Q(x, a), \forall x, a$, asymptotically over time. Our theoretical results study the behavior of $Q(x, a)$ under
57 different operators, establishing the benefits of our RSOs over other (deterministic) operators such as in [5].