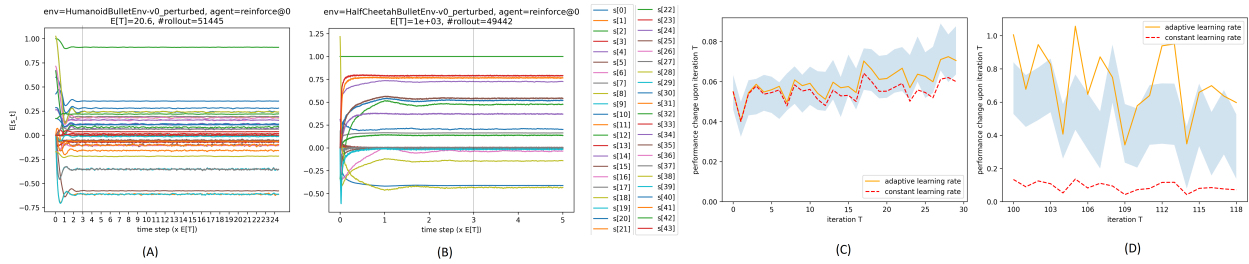


1 We thank all reviewers for the many helpful comments to our paper. A common suggestion from most reviewers is to  
 2 further elucidate how our theory helps develop better RL algorithms. We highlight two such algorithmic ideas below.

3 **Recursive perturbation.** Figure 1(e) in the appendix showed that in the Humanoid environment, the marginal distri-  
 4 butions of  $s_t$  may take as many as  $20 \cdot \mathbb{E}[T]$  steps to converge to the steady state (where  $\mathbb{E}[T]$  is the average episode  
 5 length), while Figure 1(f) showed that the marginal distributions of HalfCheetah do not converge at all (because  
 6  $\mathbb{E}[T] \equiv 1000$  in HalfCheetah, a situation where the marginal distributions will be strongly phased, as discussed in  
 7 Section 4). These observations illustrate the *necessity* to apply the  $\epsilon$ -perturbation trick in real-world RL practice.

8 In Section F.2 of the appendix we introduced an augmented variant of  $\epsilon$ -perturbation. The idea is to recursively  
 9 apply the one-shot perturbation to already perturbed models, which in the limit is equivalent to a perturbation with  
 10 self-looping probability  $\epsilon$ . We empirically found that recursive perturbation with  $\epsilon = 1 - 1/\mathbb{E}[T]$  appears to force  
 11 convergence of the marginal distribution at  $t = 3\mathbb{E}[T]$  in *various* environments. Figure 2(b) in the appendix already  
 12 showed such a result on a synthetic but challenging environment. Figure (A) and (B) below further demonstrate the  
 13 same observation in Humanoid and HalfCheetah. Comparing Figure (A) with Figure 1(e), and Figure (B) with Figure  
 14 1(f), we can see clearly the effectiveness of the recursive perturbation on these two popular RL environments.



15

16 **Episode-length-adaptive policy gradient.** Our steady-state policy gradient theorem (i.e. Theorem 5) showed a pro-  
 17 portionality factor of  $\mathbb{E}_\pi[T] - 1$  between the true policy gradient  $\nabla J_{epi}$  and the classic policy gradient estimator  
 18  $F(\theta) = \mathbb{E}_{s,a \sim \pi} [Q_\pi(s,a) \nabla \log \pi(s,a;\theta) | s \notin \mathcal{S}_\perp]$ , and we have argued in Section 5 (line 320-326) that this pro-  
 19 portionality factor can change dramatically in practical RL training, and that while it will not change the gradient  
 20 direction, using  $F(\theta)$  to estimate the policy gradient is equivalent to applying a learning rate  $\beta = \frac{\alpha}{\mathbb{E}_\pi[T] - 1}$  to the truly  
 21 unbiased estimator, where  $\beta$  is dynamically drifting along with the changes of episode length during the training.

22 Figure (C) and (D) above revealed how quickly the drifting episode length can hurt the quality of gradient estimations  
 23 in the Hopper Environment. As with Figure 3 in appendix, the shaded area gives 90% confidence intervals of the true  
 24 changes of policy performance in iterations with  $\alpha = 5 \times 10^{-4}$ . The red dotted curve is the predicted performance  
 25 change by assuming the proportionality factor as a constant absorbed in the learning rate, which quickly leads to bias  
 26 after only tens of iterations, as Figure C shows. The bias becomes quite significant after 100 iterations, as Figure  
 27 (D) shows. On the other hand, the orange curves are the estimated performance change by following the unbiased  
 28 estimator given by Theorem 5, which leads to much less bias (see Section F.3 for more details about this experiment).

29 Meanwhile, we try to give in this response a better **summary of the theoretical implications** of our paper:

- 30 1. **(a)** We identified two formal properties (Definition 1) in the learning environments of finite-horizon decision tasks.  
 31 **(b)** We proved that these properties imply existence of nondegenerate stationary distribution *in finite-horizon tasks*.
- 32 2. **(a)** We proposed a perturbation trick (one-shot version in Definition 2, recursive version in Section F.2) and proved  
 33 that the perturbed model (of learning environment of finite-horizon task) is ergodic. **(b)** The *guaranteed* (instead  
 34 of assumed) ergodicity implies that few-step sampling is sound in *the perturbed model* of all finite-horizon tasks.
- 35 3. **(a)** We analyzed the Bellman equation under the steady state, which connects, for *any* function, its mean values over  
 36 the state space and over the episode space. **(b)** As two special cases of 3(a), we showed  $J_{epi}(\pi) = J_{avg}(\pi) \cdot \mathbb{E}_\pi[T]$   
 37 and  $\rho_\pi = \mu_{epi}^1$ , which in turn help unify RL formulations between continual and episodic tasks (Table 1).
- 38 4. **(a)** We analyzed the differentiated Bellman equation under the steady state, which leads to  $\nabla J_{epi} = (\mathbb{E}_\theta[T] - 1) \cdot$   
 39  $\mathbb{E}_{s,a \sim \rho_\theta} [\dots]$ . **(b)** 4(a) implies two sources of bias in existing policy gradient algorithms as discussed above, i.e. the  
 40 fluctuating marginal distribution in tasks like HalfCheetah and the drifting proportionality factor.

41 *To Reviewer 3:* We fully agree that continual and episodic RL are equally important, yet please note that the eight  
 42 results of this particular paper (see the list above) apply only to episodic RL with finite horizon. Although we did use  
 43 infinite-horizon MDP as an *approach*, the *problems* we solved seem to have limited overlap with infinite-horizon RL.

44 *To Reviewer 4:* We hope the summary above better explained what the paper accomplished. The positioning of the  
 45 paper is actually general, with 6 out of the 8 result items applicable to both policy-based and value-based methods  
 46 (e.g. parallel few-step sampling was used in both methods [14]). Please see Section B for more discussion on this.