

1 We would like to thank the reviewers for their time. We want to reiterate our primary objective is to **make PBT more**  
 2 **efficient**, such that it can be **effective with a small computational budget**. This was motivated by our own frustration  
 3 with brittle RL hyperparameters, without the budget to run PBT. It is great to see all reviewers got this: R1 and R2  
 4 both saw the benefits in RL, saying this "might be of particularly importance for reinforcement learning algorithms",  
 5 and "empirical results shown in the RL setting are very good" respectively. Meanwhile R2 noted the improvement in  
 6 particular vs. PBT, and R4 appreciated both the theoretical results and variety of experiments. We now seek to clarify  
 7 all concerns, and hope this is sufficient to consider raising scores.

8 **BO Comparison R1 R2** We tested vanilla BO (sequential, EI, RatQuad) on the  
 9 BipedalWalker task. We used 500k samples for each trial, to have sufficient  
 10 number of evaluations in 4M and 8M timesteps (to compare vs.  $B \in \{4, 8\}$ ). In  
 11 Fig 1 we see BO performs poorly with the 4M budget, and for the 8M budget still  
 12 underperforms PB2. This sequential method should be superior in performance  
 13 to batch BO methods, since each trial is completed with full knowledge of  
 14 all previous trials. The improvement therefore represents gains from updating  
 15 hyperparameters on the fly (also see Fig 5 d) of Jaderberg et al. 2017). We also  
 16 tried to compare against BOHB, however the ray tune version does not work for  
 17 RL. We do not think it is fair to compare codebases for RL (see: Engstrom, ICLR 2020), so will seek to resolve this in  
 18 time for the CRC. In the ASHA paper they outperform BOHB, so we feel this is a strong non-PBT baseline.

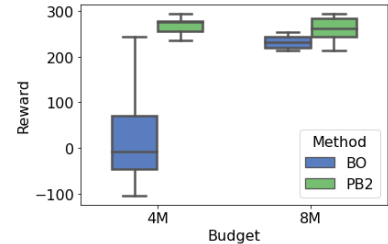


Figure 1: BipedalWalker, best rewards.

19 **Larger Population Sizes R2 R4** We agree it is important to assess PB2 with  $B > 8$ . To  
 20 answer this, we ran the BipedalWalker experiment with  $B = 16$ . As we see in Figure 2,  
 21 both PBT and PB2 achieve optimal rewards ( $> 300$ ), but PB2 is still more efficient. We  
 22 hope this provides confidence in our method's effectiveness with more resources.

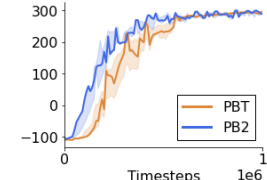


Figure 2: BipedalWalker, median, 5 seeds,  $B = 16$ .

23 **Larger Supervised Experiments R2 R3** We included the SL experiment to show our  
 24 method generalizes beyond RL, as noted by R4. The relative performance of the hyper-  
 25 parameter algorithms should be agnostic to the architecture, and PB2 outperforms other  
 26 competitive baselines **across five seeds**. We used a medium sized CNN, as this was the  
 27 most powerful model we could realistically train with our compute. Larger networks are prohibitively expensive for  
 28 our lab. Instead, we allocated our resources on larger RL settings such as Atari (which are more CPU intensive). For  
 29 reference, a similarly accurate CNN model was used for the first set of experiments in ASHA (as a POC). The main  
 30 result in BOHB was a CNN on CIFAR-10, but they used 19 workers, each with 2 GPUs. However, their RL experiments  
 31 were toy. In light of this, we think expecting SOTA SL experiments on a paper focused on RL is unreasonable.

32 **What is  $fg_{t+1}(x)$ ? R1 R2** It should be written as  $g_{t+1}(x)$ , without  $f$ . Thank you for catching this.

33 Next we address individual comments in more detail.

34 **R1: Wall-clock comparison** For RL, the GP-bandit step takes a trivial amount of time compared to querying a simulator  
 35 or computing gradients, given we only have  $< 10$  hyperparameters we have a tiny dataset. If we have vast computational  
 36 resources and we wish to optimize  $\gg 10$  hyperparameters for  $B > 20$ , it may be required to use more scalable GP  
 37 methods. However, *this is not the problem we are trying to solve*. **PB2 only slightly outperforming ASHA** The main  
 38 purpose of this paper is to improve PBT, and our main goal is to show this. In addition, we also show PB2 outperforms  
 39 state of the art (ASHA). **Sec 3.1** To be precise, if (1)  $\omega = 0$  and (2) the number of parallel agent  $B = 1$ , then our model  
 40 reduces to GP-UCB. The reduction factor is 3, this is the default in the ray tune implementation.

41 **R2:** Thank you very much for your comments, we are glad you feel this is an improvement over PBT, in particular  
 42 making it more usable, which was always our goal. We dont assume the convexity of the function. Instead, we make a  
 43 mild condition on the compact and convex input space so that the input space is continuous and not segmented, e.g.,  
 44 the case of non-convexity. This mild assumption is popular in GP bandit literature (see Srinivas et al, 2010). We  
 45 always have  $B \ll T$  because the number of update step  $T$  typically goes beyond thousands to millions while the  
 46 batch size  $B = 4, 8, 16...$  is much smaller and limited to the parallel facility we have. **Hyperparameter Schedules** We  
 47 disagree the hyperparameters are random, for the learning rate in particular, you can see the GP balances exploration  
 48 and exploitation by either selecting points from a single mode or right at the boundaries.

49 **R3: Distinguishing features of our proof** Our proof makes two significant extensions beyond TV-GP-UCB (Bugonovic  
 50 et al, 2016). First, we improve the bound from C.2 from  $\tilde{N}^3$  to  $\tilde{N}^{2.5}$  (smaller is tighter and better). Second, we extend the  
 51 regret bound to the batch setting including the new Lemma 8,9,10 (in the Appendix). **Section 5.3** The hyperparameters  
 52 and ranges are shown in Table 9, in the Appendix (top of p14). **Larger datasets using larger networks** See above: we  
 53 hope this is not the primary reason for rejecting our work. **Optimal network architectures**, given recent works using  
 54 BO for NAS we think this is a logical next step, and something we are excited about. NAS for RL is a nascent field.