

1 We thank anonymous reviewers for their valuable comments and suggestions.

2 **Comment 1: Training time (Reviewer #1)**

3 **Response:** We have included the training time in the paper (L318) that “Within one million timesteps, the training
4 wall-clock time for our TRGPPO is 33 min; for PPO, 32 min”. We use several techniques to allow efficient optimization,
5 including the problem reduction, the DNN-approximation and problem discretization (described in Sec 5.1).

6 **Comment 2: Concerns of the experiment evaluation, random number and baseline (Reviewer #1 & #2 & #3)**

7 **Response:** We used the averaged top 10 reward following the setting in [24], which could somewhat reflect the
8 algorithm’s ability on searching good solution but we agree it’s somehow unreliable. However, we have also plotted
9 the learning curves in Fig. 3, which could help infer the stability of the algorithm. Per your suggestion, we have
10 made several revisions, listed as follows: 1) *report the averaged reward over all episodes of training*; 2) *compare with*
11 *the baseline of adaptive KL regularization of PPO* (and clarify the related description in the introduction); 3) *run a*
12 *hyperparameter sweep for ϵ of PPO over $[0.1, 0.6]$ with step 0.05*; 4) *increase the number of random seeds to 10*. We
13 normalized the scores for each environment so that the random policy gave a score of 0 and the best score was set to 1.
14 The averaged normalized scores (over 60 runs with all episodes of training for each algorithm, on 6 environments) are
15 as follows: **TRGPPO: 0.629**; *PPO($\epsilon = 0.2, \text{default}$): 0.441*; *PPO($\epsilon = 0.25, \text{optimal PPO}$): 0.484*; *PPO-adaptiveKL:*
16 *0.422*. We will add more details of the results in the final version.

17 **Comment 3: Concerns about Lemma 2 (L113) and several typos (Reviewer #1)**

18 **Response:** Thanks for your careful reading. The correct form of the LHS of the equation in Lemma 2 should be
19 $\mathbb{E}_{\pi_{t+1}} [\pi_{t+1}(a) | \pi_0]$. This typo does not affect the correctness of the lemma and the remaining theoretical results in the
20 manuscript. We will rectify all the typos in the final version.

21 **Comment 4: The existence of $\pi_{\text{new}}^{\text{PPO}}$ (L235) (Reviewer #1)**

22 **Response:** The problem is that how we can find $\pi_{\text{new}}^{\text{PPO}} \in \Pi_{\text{new}}^{\text{PPO}}$ that achieves minimum KL divergence on all
23 states s_t , which can be formalized as $\min_{\pi \in \Pi_{\text{new}}^{\text{PPO}}} (D_{\text{KL}}^{s_1}(\pi_{\text{old}}, \pi), \dots, D_{\text{KL}}^{s_T}(\pi_{\text{old}}, \pi))$. Note that $\pi(\cdot | s_t)$ is a conditional
24 probability and theoretically the optimal solution on different states are independent from each other. Thus the problem
25 can be optimized by independently solving $\min_{\pi(\cdot | s_t) \in \{\pi(\cdot | s_t) : \pi \in \Pi_{\text{new}}^{\text{PPO}}\}} D_{\text{KL}}(\pi_{\text{old}}(\cdot | s_t), \pi(\cdot | s_t))$ for each s_t . The final
26 $\pi_{\text{new}}^{\text{PPO}}$ is obtained by integrating these independent optimal solutions $\pi_{\text{new}}^{\text{PPO}}(\cdot | s_t)$ on different state s_t . We have provided
27 detail in Appendix D and we will add more explanation in the final version.

28 **Comment 5: Concerns about the details of and the reproducibility of the experiment (Reviewer #1)**

29 **Response:** We used Gaussian and Gibbs policy for continuous and discrete tasks respectively, parametrized by a DNN.
30 For baseline, we used the setting recommended in the original paper. We have also submitted a link of the source code
31 as supplementary (L48 in the paper). We will add more details and release our code in the final version.

32 **Comment 6: How is Eq. (4) transformed into Eq. (5) in supplementary? (Reviewer #2)**

33 **Response:** To be brief, *let’s number the equations in Eq. (4) by (a)-(d)*. First, by (a)(b), we have $\lambda \neq 0$, since if
34 $\lambda = 0$ then $\nu = 0$ (by (a)), which contradicts (b). Second, by (c) and $\lambda \neq 0$, we have $\sum_{a \in \mathcal{A}} p'_a \log(p'_a/p_a) = \delta$.
35 Third, taking (a) into (d), we have $p'_a/p_a = \nu/\lambda = (1 - p'_{a_t})/(1 - p_{a_t})$ for $a \neq a_t$. Then, taking this equation into
36 $\sum_{a \in \mathcal{A}} p'_a \log(p'_a/p_a) = \delta$, we obtain Eq. (5). We will add more details in the final version.

37 **Comment 7: The performance on Humanoid (Reviewer #2)**

38 **Response:** One possible explanation is that the larger clipping range of our TRGPPO may make it suffer from the
39 noisy estimated advantage values, especially at the later training phase where the advantage values are large and noisy.
40 This issue could be addressed using our adaptive clipping scheme by taking the trade-off between exploration and
41 stability into account. In particular, in the revised version, we have implemented two variants of TRGPPO: linearly
42 decaying ϵ from 0.2 to 0.1 (named by *TRGPPO-decay*) or clipping the clipping ranges (named by *TRGPPO-clipping*),
43 i.e., $l_{s,a}^{\delta, \epsilon, \epsilon_t} = \text{clip}(l_{s,a}^{\delta}, \epsilon_t, \epsilon)$, $u_{s,a}^{\delta, \epsilon, \epsilon_t} = \text{clip}(u_{s,a}^{\delta}, 1/\epsilon, 1/\epsilon_t)$, where $0 < \epsilon < 1$ and $\epsilon_t = \epsilon t/T$ are parameters to control
44 the level of the clipping ranges, t and T are the current and total training iterations respectively. Both these two
45 methods could improve the reward and sample efficiency. The averaged episode rewards over all episodes of training on
46 Humanoid are as follows: *TRGPPO-decay: 3013.3*; **TRGPPO-clipping: 3148.1**; *PPO: 2944.2*. The timesteps ($\times 10^3$)
47 to hit the threshold are as follows: *TRGPPO-decay: 7514*; **TRGPPO-clipping: 7132**; *PPO: 9088*.

48 **Comment 8: Prove that TRGPPO converges to the optimal policy (Reviewer #2)**

49 **Response:** Thanks for your suggestion. It’s interesting to prove such convergence property. However, there seems
50 does not exist closed-form of our clipping range in Eq. (5), making it hard to measure the improvement of $\Delta_{\pi_0, t}^{\text{TRGPPO}}$
51 by the explicit form of $\mathbb{E}_{\pi_{t+1}} [\pi_{t+1}(a_{\text{opt}}) | \pi_t]$ (see Eq. (3)). Alternatively, we plan to work on this by analyzing the
52 corresponding bound for each term in Eq.(3).

53 **Comment 9: Some mathematical formulae in the paper could be better formatted. (Reviewer #3)**

54 **Response:** Thanks for your comment. We will polish mathematic notions and align expressions in the final version.