

1 We thank the reviewers for the reviews, providing meaningful insight with constructive feedback. Due to the page
2 limitation, we only provide critical values. In the final manuscript, we will present all the results, alongside modifications
3 reflecting reviewers' comments which are not mentioned in this response.

4 **R1: Result on Humanoid environment.** We tested our method on Humanoid-v2 and confirmed our method works
5 properly. The relative baseline and adaptive variance update algorithm, which performs best among proposed algorithms,
6 was tested on the environment and scored 6169 ± 455 with ten random seeds at 1M steps.

7 **R1: Interaction of ES and PG workers.** We measured how many RL and EA actors were contributed in improving
8 the performance, as a summation of the update ratio p (Eq. 6), with higher value indicating more contribution. In our
9 method, RL actors contributed twice more compared to ES actors in HalfCheetah, with values of **214.53** and **105.52**,
10 respectively. The result was reversed in Hopper, where RL contributed **200.86** while EA actors did **363.53**.

11 **R2, R3: Evaluation method for performance and speed.** We evaluated our algorithm in two perspectives; perfor-
12 mance improvement and speed improvement. For the performance improvement, we evaluated our method as same as
13 the baselines for a fair comparison. In many papers, the final score of the fixed interaction step is frequently used for
14 evaluation metrics. Therefore, all performance result scores are measured in the fixed interaction step. For the speed
15 improvement, we measured execution wall-clock time for the fixed interaction step; the result is presented in Table 1.
16 As the CEM-RL is implemented in a serial-synchronous, we modified the algorithm to a parallel-synchronous version
17 (P-CEM-RL), and then we compared these algorithms with our method to show the efficiency of the asynchronism. We
18 will include the learning curve in the final version.

19 **R2: Ablation study is missing.** Our algorithm mainly consists of three aspect; asynchronism, mean and variance
20 update rules. We presented the effect of a simple asynchronous method [25] with the CEM-RL update rule in column
21 "Rank-based" of Table 2. We presented the effect of the variance update rule in Appendix C.3 by comparing the result
22 with a fixed variance setting. Then, we provided all combinations of our proposed mean and variance in Table 2.
23 However, we agree with the reviewer because all these results are shown separately and not discussed thoroughly in the
24 manuscript. We will add a section so that it can be seen at a glance. If these results are still not enough for an ablation
25 study, it would be beneficial for us to consolidate our manuscript if the reviewer can provide us more specific guideline.

26 **R2: Asynchronism in CEM-RL is not "impossible".** We used the word "impossible" to emphasize that the update
27 rule of CEM-RL cannot be used exactly the same in an asynchronous setting. CEM-RL spawns all actors at the same
28 time with a fixed number of each agent. In an asynchronous setting, some modifications should be applied, such as
29 alternatively spawning RL and ES actors. We intended to highlight the difference; however, we agree that the word is
30 too aggressive. We will soften the tone in the final version.

31 **R2: (1+1)-ES is not a fitness-based method.** We categorized (1+1)-ES as a fitness-based method because it uses
32 fitness values for comparing. However, as the reviewer pointed out, the method is ambiguous to be categorized. We will
33 add detailed explanation in the final version.

34 **R2: No reference or evidence about the statement "Capability of high exploration in fitness-based scheme".**
35 Our phrase "high exploration" is intended to empathize with the aggressiveness when there appears a superior individual.
36 However, there might be a misunderstanding about the general meaning of exploration in the policy search field. We
37 will modify the statement in the final version.

38 **R3, R4: Speed up is the factor of 2 or 3. Time efficiency experiment was conducted with only one setting.** We
39 used five actors as provided in the section 4.1. The number of actors is not limited, but it requires GPU calculation,
40 limiting the experiment in our setting. We additionally tested our methods with various actors of 2, 3, 4, 5, and 7
41 in Halfcheetah. The running times were **75%**, **42%**, **37%**, **32%**, and **25%** compared to the execution time of the
42 CEM-RL. It shows that time efficiency is linearly increased as number of actors increases. We will provide broader
43 experiments and discussions about the number of actors, with a table and also a graph.

44 **R3, R4: Ask and Tell based update rules are missing. Using variance instead of covariance.** We took a look into
45 the Nevergrad library and read the original papers of implemented algorithms. We will try to merge various algorithms
46 that fit the update scheme of combining ES and RL. It seems though some algorithms are hard to be merged. Our
47 network consists of only three layers with 100k parameters, which is very shallow compared to the networks in the
48 computer vision field. However, algorithms that use covariance like CMA-ES are not appropriate with 100k parameters.
49 It is also the reason why the baseline CEM-RL used variance instead of covariance.

50 **R4: Stability metric is not provided.** We defined the term "stability" for consistently showing high performances,
51 thereby having low std value per mean (σ/μ). Our final method AES-RL, relative baseline with the adaptive update,
52 is chosen because its performance is high, and the σ/μ value is low. As shown in Table 2, previous algorithms for
53 asynchronous updates show a higher σ/μ compared to our method, therefore we claimed that our method is relatively
54 stable. We will explicitly explain the metric in the plain text and also emphasize it in the tables.