

1 We thank all the reviewers for the time and expertise invested in these reviews. We have corrected the typos and
2 grammar mistakes, and rephrased some points/sentences to improve the overall readability of our paper. In the following
3 sections, we respond to the major concerns of each reviewer respectively to, hopefully, clarify our methods and visions.

4 **Responses to Review # 1**

5 **Q: What is the meaning of every notation?** A: We are sorry that some abuse of notations in the paper hinders the
6 understanding of our method. We have checked the notations to ensure that they are consistent throughout the whole
7 paper. We would further explain that, generally, we use K to denote the number of ensemble models and H to denote
8 the maximum horizon. Their corresponding lowercase letter refer to one instance in the set, e.g. s_{t+h} means a state
9 in $\{s_{t+1}, \dots, s_{t+H}\}$. Specially, s_g means a “good” state in s_{t+1}, \dots, s_{t+H} that can be regarded as a goal state (to be
10 consistent with notations in other goal-conditioned work).

11 **Q: What is the relationship to other Transfer Learning/Imitation Learning method?** A: This work aims to tackle a problem
12 that lies in the intersection of Imitation Learning and Transfer Learning (more specifically Sim2Real). We have included
13 the most relevant works that tackle a similar problem and demonstrated our difference and novelty. We would like to
14 conduct a complete literature review later to cover the recent works in Imitation Learning and Transfer Learning and
15 add them to the related work section.

16 Since there are no major flaws pointed out in the review, could the reviewer please raise the overall score?

17 **Responses to Review # 2**

18 **Q: Can this method work in real-time control requiring a reset-able simulator?** A: Generally, the computation cost of a
19 reset-able simulator is comparable to a model-based method and is thus acceptable. We will explore how to relax such a
20 constraint in future work from both methodological and engineering perspectives.

21 **Q: Can this method work considering the complexity of real dynamics mismatch?** A: Our empirical results on various
22 modifications to MuJoCo environments (3 types \times 3 magnitude) can prove that our method is robust to different
23 dynamics mismatch, so we believe it shows the potential to work in sophisticated real-world problems. We have plans
24 to apply the method proposed in this work on a real quadruped robot in the future.

25 **Responses to Review # 3**

26 **Q: What is the difference compared with others using Goal-conditioned Policy(GCP)/Hindsight Inverse Dynamics(HID)?**
27 A: HID adopted in our work contributes to the overall purpose to alleviate dynamics mismatch problem and augment the
28 limited expert data. Other works use GCP/HID for different purposes. PCHID (arXiv 1910.14055) solves goal-oriented
29 tasks in a supervised manner. Play-GCBC (arXiv 1903.01973) trains a goal-conditioned policy to address a multi-modal
30 problem. Relay Policy Learning (RPL) mentioned in (arXiv 1910.11956) solves long-term robotic tasks in a hierarchical
31 manner with the help of a goal-conditioned policy. We will discuss these studies in the related work section.

32 **Q: Can this method work considering the complexity of real dynamics mismatch?** A: Our empirical results on various
33 modifications to MuJoCo environments (3 types \times 3 magnitude) can prove that our method is robust to different
34 dynamics mismatch. So we believe it shows the potential to work in sophisticated real-world problems. We have plans
35 to apply the method proposed in this work on a real quadruped robot in the future.

36 **Responses to Review # 6**

37 **Q: What is the meaning of “partial alignment”** A: We realize that the word “partially” in Line 9 is a little bit confusing
38 and have clarified it. Instead of referring to environment being “partially observable”, it means that not every (s_t, s_{t+h})
39 pair needs to be matched, and a subset (which “partial” actually means) matching is enough. Please refer to Figure 3 in
40 the paper for a visual illustration. All the experiments are done in fully observable MuJoCo environments.

41 **Q: What is the rationality behind “partial alignment” assumption?** A: Such an assumption comes from an empirical
42 observation that in robotics control problems, some key poses in different dynamics are still alike. Aligning such
43 key poses would make the long-horizon learning much easier. Besides, a weighted GAIL(GAILfO) optimization in
44 our method is a relaxation to the “exact partial alignment”. It can be viewed as an occupancy matching problem that
45 automatically matches the most “similar” state. We would like to explore further in this area to see if this method will
46 work under some formalization of dynamics differences.

47 **Q: Do we use the simulator reward in the training phase?** A: No, we do not use simulator reward functions in any training
48 phases, i.e., both GAILfO and Goal-conditioned BC (which accords with the usual Imitation Learning setting).