1  We thank the reviewers for lending their expertise and time to provide feedback on our efforts. We are glad that all the
2  reviewers found our insight that action transformation can be seen as an IfO problem novel and interesting. We respond
3  to the biggest questions and comments below and will address all feedback in the paper.

4  [R2 , R3 , R4 ] The reviewers are correct in pointing out that, despite the title, we do not include a real robot experiment.
5  Our work is motivated by sim-to-real, but we were unable to conduct real robot experiments due to the current pandemic
6  as R1 and R4 pointed out. If accepted, we will make several changes to moderate the claims as R4 suggested. In
7  particular, we will change the terminology in the paper to align with more general transfer learning, using source
8  and target domains as opposed to sim and "real." Also, we will change the title to "Towards Sim-to-Real Transfer:
9  An Imitation from Observation Approach." Please note that our formulation remains very relevant to the sim-to-real
10  community. We would like to highlight that one of our experiments is indeed an excellent proxy for the sim-to-real
11  problem: In the Minitaur domain (Figure 2), Tan et al. [38] found that while their existing simulator (our source domain)
12  inaccurately represented their robot, the new simulator they crafted (our target domain) *did* enable direct policy transfer
13  from sim to real.

14  [R2 , R3 ] Both manipulation domains [7, 24, 26, 39, 40, R2 's suggestions] and locomotion domains [9, 10, 11, 14,
15  18, 27, 38, 46] are prevalent in the sim-to-real literature. Both are important—but different—problems: manipulation
16  domains are more likely to exhibit observation mismatch, whereas locomotion domains are more typically associated
17  with dynamics mismatch. The scope of our work here is mainly dynamics mismatch, and therefore we focus our
18  experiments on locomotion problems. GARAT solely addresses dynamics mismatch. For locomotion, the observations
19  are usually joint angles and velocities, so observation mismatch is negligible. If accepted, we will make this scope more
20  clear in the camera-ready version of our paper and include the references R2 suggested. Note that in our problem setting,
21  the state spaces are the same in the source and target domains, as is commonly the case in sim-to-real. Specifically, we
22  consider dealing with embodiment mismatch to be beyond the scope of this paper.

23  [R2 ] Most domain randomization techniques, and all the papers suggested by R2 , require a modifiable simulator and
24  substantial domain expertise [7]. In this paper we focus on the case where the simulator cannot be modified (black box),
25  and hence it is not appropriate to compare with methods that can adjust the simulator itself. We compare to ANE [20]
26  which is an action randomization technique.

27  [R2 ] Respectfully, we strongly disagree with the reviewer's assertion that our approach does not offer significant
28  technical novelty. In this work, we show how tools developed in the imitation learning community can be successfully
29  adapted to sim-to-real problems. Moreover, our adaptation of one such tool actually leads to better performance
30  than alternative applicable approaches. To the best of our knowledge, this is the first time this has been studied
31  in the literature, and therefore our work represents a novel and important connection between two largely separate
32  communities.

33  [R1 ] Concerning why GAT was not as effective as GARAT on transfer, perhaps it would be useful to compare the two
34  techniques to their imitation learning equivalents, behavioral cloning (BC) and using inverse RL (IRL). BC suffers
35  from distribution shift while IRL methods are able to learn how to recover from such shifts; likewise, GAT is unable to
36  recover from the shift introduced by an imperfect action transformation while GARAT can correct for such deviations.
37  GAT is myopic, trying to match single transitions, while GARAT matches the whole trajectory (Figure 1b).

38  [R2 , R3 ] The curve for GAT cuts off early in Figure 1b. In the InvertedPendulum domain, the episode terminates if
39  the angle of the pendulum exceeds $\pm 0.2$ radians. In the environment with GAT, the action transformation learns to
40  keep close to the target domain's dynamics early on, but this causes instability later in the episode, leading to early
41  termination. GARAT sacrifices initial accuracy to keep the overall trajectory as realistic as possible. We will edit the
42  caption for Figure 1b to make it clear in the camera ready version of the paper.

43  [R3 ] We use the loss derived in Section 4.3 in our main results. Our algorithm is agnostic to the RL algorithms used
44  for training. We chose PPO and TRPO for the action transformation function and the agent respectively because that
45  combination worked best in preliminary experiments on the InvertedPendulum domain.

46  [R3 ] GARAT should implicitly address process noise due to its adversarial learning procedure. The discriminator in
47  GARAT encourages the action transformation function to learn a distribution of transitions that are similar to the target
48  domain, including any noisy transitions. Moreover, GAT [14] has been shown to be useful in sim-to-real transfer on a
49  real legged humanoid robot, showing that impact dynamics and operational noise do not prevent learning.

50  [R2 , R3 ] Figure 3 was normalized in order to compare the performance of different algorithms across different
51  domains. It does not represent the maximum and minimum returns possible. We train $\pi_{real}$ in the target domain for 1
52  million time-steps, enough to reach a reasonable policy. These policies may take more training to converge completely
53  (HalfCheetah is usually trained for 10 million timesteps). GARAT manages to learn a policy that does better than the
54  policies trained directly in the target domain for some of these environments.