1  We thank the reviewers for their feedback. We are excited to see that all four reviewers have rated our work positively
2  and found that our paper addresses a very important problem (R1), is nicely motivated and builds logically on prior
3  work (R2), works significantly better than baselines (R3), and is clearly contrasted with existing work (R4). We address
4  the most salient points of feedback below, and will incorporate all of their feedback in the final version of our paper.

5  ——————————————————————— **General Feedback** ———————————————————————

6  **Novelty** (R1, R2). "R2: *similar methods ... were published recently. It is very close to ... HER and the presented ideas*
7  *for sampling new reward parameters are straightforward.*" Indeed, the power of our work is that we present simple
8  methods to extend HER to general multi-task problems. Work by Eysenbach et al. 2020 was concurrent with ours, and
9  is likely simultaneously in submission to NeurIPS 2020. Nevertheless, we contrast our work with HIPI as follows: our
10 work explores different environments, has a different connection to inverse RL, and has extensive comparison of the
11 pros and cons of each method. We're running HIPI as a baseline for all of our environments, which we will add when
12 complete. However, Fig. 6 shows HIPI is worse on easy environments while taking more time to train; hence it is hard
13 to see it working for more complex settings. Computationally, HIPI is $\approx 11\times$ slower than our method on environments
14 that require bigger batch sizes. On Ant and Humanoid, this requires almost two weeks to run, which is limiting in
15 academic labs such as ours. This analysis of computational complexity will be added to the paper.

16 **Tasks** (R3, R4). "R4: *the authors mostly use customized environments with dense rewards and goal-conditioned tasks*
17 *where the task is represented by more than the goal state.*" For sparse goal-reaching, HER should work better, since
18 the optimal relabeled task can be easily computed from the states reached. Our algorithms are designed for settings
19 beyond pure goal reaching, where HER cannot be straightforwardly applied. We will clarify the categories of rewards
20 that go into each environment, beyond the explicit reward formulae in Appendix B. And although our methods may not
21 accelerate learning two highly dissimilar tasks, e.g. hammering a nail and opening a door, we show that they work well
22 on a variety of meaningful and challenging robotics problems, such as control with safety or energy trade-offs.

23 **Theory** (R1). "R1: *There are no bounds that show how well both algorithms approximate the IRL problem and it's very*
24 *unclear which one is more accurate.*" As with most prior work in hindsight relabeling [2, 46, 38, 45], we found it difficult
25 to prove any meaningful bounds, likely due to the complexity of handling a distribution of policies $\{\pi(\cdot|z)|z \in \mathcal{Z}\}$ and
26 the accuracy of the $Q$-function. Both are complicated by function approximation and complex dynamics and rewards.
27 Intuitively, AIR may be more accurate because AIR directly maximizes the max-margin IRL objective when comparing
28 against infinite random trajectories in the limit. We compared our relabeling algorithms empirically in Sec 4.4, and
29 would love to see future work advancing our theoretical understanding of hindsight algorithms.

30 ——————————————————— **Algorithmic / Experimental details** ———————————————————

31 **Negative trajectories** (R1). "R1: *the authors' algorithm would always pick the successful trajectories even though we*
32 *know that informative negatives are crucial for off-policy RL algorithms.*" Indeed, hindsight relabeling in stochastic
33 environments introduces hindsight bias, a tradeoff we discuss in Appendix C. In spite of this, our relabeling methods
34 achieve good performance in the robotics environments that we test in. Extending our methods to more stochastic
35 environments, e.g. by applying ARCHER [37], is an exciting avenue for future work.

36 **Relabeling trajectories vs transitions** (R1). "R1: *It is not clear why relabeling is done on the entire trajectory rather*
37 *than on the individual transition.*" It's possible to relabel single transitions, based on (a) the transition reward or (b) the
38 $Q$-value, but these both have problems. (a) Relabeling just based on the 1-step reward produces myopic relabeling that
39 values short term rewards over large rewards that accumulate later. (b) In Sec 4.5, we present an in-depth discussion of
40 the problems of using the $Q$-value to do transition relabeling. Relabeling the entire trajectory once is computationally
41 efficient, simple, and already provides large performance gains. We will add this discussion earlier in the paper.

42 **Distribution shift** (R3). "R3: *The trajectories may be generated from a different policy parameter from that one*
43 *being optimized.*" This is mostly fine for off-policy learning algorithms, such as SAC. We discuss bootstrap error in
44 Appendix C, where we note that distribution shift from relabeled trajectories should at worst result in overestimation
45 and encourage exploration of promising areas of the MDP. Empirically, our methods still perform well, and a promising
46 extension of our work is to use offline RL methods to further reduce distribution shift error.

47 **Hyperparameters for relabeling** (R4). "R4: *There are some hyperparameters (e.g. number of tasks returned $m$) that*
48 *need to be specified when deploying AIR.*" We simply select $m = 1$ relabeled trajectories out of $K = 100$ sampled task
49 variables for all environments, and will add an ablation study on the effect of $m$ and $K$ in the camera-ready's appendix.
50 With $m = 1$, our relabeling algorithms are empirically accurate in finding the proper relabeled task, as seen in Sec 4.4.

51 **Error bars** (R4). "R4: *I didn't notice anything specifying what this error bar signifies.*" The error bars show the
52 standard deviation across seeds. When corrected by $1/\sqrt{n}$ for the number of runs, they have minimal overlap.

53 We again thank the reviewers for their detailed reviews, and even pointing out a few typos and missing citations, which
54 we will add to the final version of the paper.