

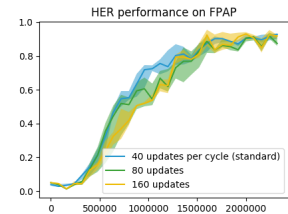
1 Firstly we would like to thank all of the reviewers for their insightful and thorough feedback. We really appreciate the
2 time taken to go through the paper and will try our best here to respond to as many of the questions raised as possible.

3 To start with, three of the reviewers noted that we had not included any comparisons with other model-based approaches.
4 Our initial reasoning for not including such a comparison was that because standard model-based methods generally
5 do not outperform model-free methods in terms of their final performance, it seems implausible that a model-based
6 approach designed for standard environments (i.e. dense reward, not goal-based) would perform better than, e.g., DDPG,
7 which we do include as a comparison and which performs badly on the more challenging environments. Having thought
8 about this more, however, it does seem natural that such a baseline should still be included to verify this intuition, and
9 we will look to add such a baseline in the next draft. We also thank **reviewer 1** for drawing our attention to a number of
10 papers that we missed out in the related work.

11 **Reviewer 2** raises a point about the baselines potentially not being anchored with previous numbers. Whilst it's true that
12 we ran all of the baselines ourselves, we nevertheless used the author's publicly released code for all of the baselines,
13 except HER. For HER, we used the OpenAI Baselines implementation. If we compare with Plappert et al.¹, they
14 show results for FetchPush and FetchPickAndPlace in their Figure 3 (note 1 epoch = 95000 environment interactions).
15 Although perhaps not immediately obvious, if you carefully compare our results with theirs using this scaling you
16 should see that they (at least approximately) match up. To address some of reviewer 2's other questions: we agree that
17 it would be an interesting line of work to explore similar approaches for more general environments. The reasons we
18 focused on sparse reward, goal-conditioned environments is that it felt natural to combine the training of the GANs
19 with hindsight experience relabelling (which is really at the core of our approach), as well as the fact that typical
20 model-based approaches are likely to be inappropriate for solving these kinds of task. Regarding the comment about
21 1199, it's true that other strategies can be used here. In our initial experiments we did try simply using a single future
22 goal for the whole unrolled trajectory, but found this did not perform quite as well. At the very least this should have
23 been mentioned, and so we will include a comment and perhaps consider adding an ablation demonstrating this. We
24 appreciate the comments re: 1195, 1217 and 1225 and agree that your suggested terminology would be clearer. Re:
25 the comment about 1227, we think this is a slight misunderstanding and something we need to make clearer in the
26 text. All GANs in the ensemble are trained from batches taken from the same replay buffer. For each step of each
27 imaginary rollout we choose a random GAN in the ensemble, however these imaginary rollouts are not stored in the
28 buffer (which only stores the actual rollouts with the final actions chosen by the planner, along with some initial random
29 trajectories). This is related to a comment by **reviewer 3** who asks for clarification about this as well as asking how
30 important "mixing" GANs in this way was. Essentially, this is the primary way in which the ensemble (rather than just
31 using a single GAN) was made use of, so in that sense the ablation in Appendix B demonstrates that it does have a
32 somewhat significant effect. However, it is true that there are other ways we could have made use of the ensemble,
33 but we have not considered these so far. Re: the comment about 1235, R is a typo (thank you for spotting this). We
34 weren't 100% sure what you meant by "... (not an integer, depends on ordering...)", but just to clarify (and we will adopt
35 this notation in the next draft): the planning process gives a score, n_i to each of the Q initial seed actions a_i (note the
36 actions are continuous). We then define weights $w_i = e^{\alpha n_i}$ and return $a_t = \frac{\sum_{i=1}^Q w_i a_i}{\sum_{j=1}^Q w_j}$.

37 To address some of **reviewer 3's** other questions: the training curves do indeed account for the initial random trajectories
38 (including those we discard and don't store in the replay buffer). M in figure 1 is supposed to represent the one-step pre-
39 dictive model, although we realise now there is an error in this (M should have both s_t and a_t as inputs) — we will fix this.
40 We did consider removing the OSM regularization, but decided not to as we felt it was
41 still an interesting experiment, despite the fairly small difference it made (having said this, the difference on FPAP was $\sim 5 - 10\%$, which looks small on the plot but is not entirely
42 insignificant). Given that there was not much difference between noOSM and OSM, it's
43 also the case that there is not any strong dependence on λ .

44 Addressing some of **reviewer 4's** comments: it's true that we use more updates than HER
45 when training our model-based approach. However, previous work has shown increasing
46 $num_updates$ can actually degrade HER's performance slightly. We ran some experiments
47 on FPAP to demonstrate this (and could add these to an appendix perhaps) — see the plot on the right. It's an interesting
48 question to ask whether HER could also be improved by using an ensemble — however exactly *how* you would best
49 make use of such an ensemble of policies/value functions in that context is not entirely trivial and an interesting research
50 question in itself. We feel this is also true of employing HER with a forward model (again, a very interesting suggestion)
51 — would you train DDPG+HER within a learned world model or use some other kind of model-based search to gather
52 rollouts? How often would you update the model (or continuously update)? There are a number of considerations that,
53 whilst interesting, go beyond a very straightforward baseline to compare with, in our opinion.



¹Multi-Goal Reinforcement Learning: Challenging Robotics Environments and Request for Research, 2018