

1 Thank you to all reviewers for taking the time to give us feedback on our submission. We hope the following addresses your concerns.

2 **Reviewer 1:**

3 - We would like to clarify, the *greedy* policy (exploitation) is not the one learned by the advisor. That is the task specific policy. The  
4 exploration policy (the epsilon part) is the one learned by the advisor, which as shown in Figure 6, is not simply a general policy to  
5 solve many tasks.

6 - We will expand our explanation on when we can expect this method to be useful. Intuitively, it helps when in a given state certain  
7 actions are not appropriate regardless of task variations. For example, if the pole in pole balancing is about to fall to the right, moving  
8 the cart to the left will never help the situation.

9 - Thank you for the recommendation on other exploration baselines. We will take them into account.

10 **Reviewer 2:**

11 - We also found it surprising how little work there is in this specific area. After going over your suggestion, we believe a discussion  
12 on Simsek and Barto's work is appropriate.

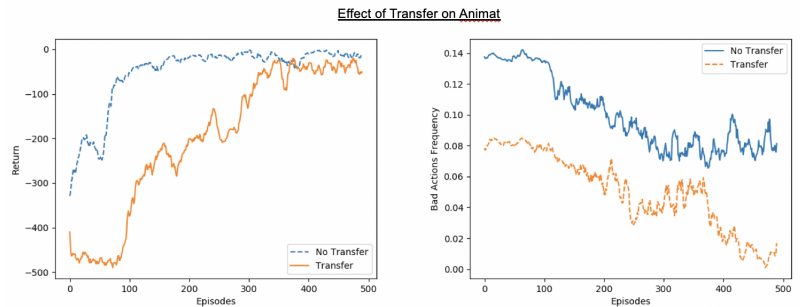
13 - We apologize if there was some confusion and we will clarify in the text. 1) Yes, we do assume that tasks come from a distribution  
14 for our theoretical formulation; however, in practice we use a finite set of tasks to solve, a small subset of which is used to learn  
15 the exploration policy. 2) No, we do not need to train in parallel. The results hold when training on different MDPs sequentially,  
16 but when possible, updating the advisor from multiple tasks in parallel makes learning an exploration policy quicker. Good point!  
17 Training on increasingly difficult problems, like in curriculum learning, would be a clear scenario where learning an exploration  
18 policy, as we propose, could lead to clear benefits.

19 - We did perform the experiment you are suggesting to evaluate the learned exploration policy. Figures 5 and table 1, evaluate the  
20 performance of the learned exploration policy on a set of novel tasks (different from the ones used for training), showing that it leads  
21 to clear performance improvement.

22 - Your comment on the Animat problem is a fair point. A simple transfer of policies would definitely reduce the  
23 number of 'non useful' actions taken. However, it would be highly biased to the task where the policy was learned,  
24 which could lead to really poor performance in a new task. We evaluated this point based on your comment.

25 The figure on the right shows the (average)  
26 effect of a simple transfer of policy over five  
27 task variations. The right plot shows that the  
28 frequency with which the agent takes 'non  
29 useful' actions is decreased significantly  
30 with simple policy transfer; however, left  
31 plot shows that a simple transfer of policy  
32 can actually make learning more difficult.

33 - We will gladly add your suggestions to  
34 improve presentation. The reason why no  
35 training progress is shown for Animat is  
36 that the plot would not show anything that  
37 was not shown for pole-balancing, and we  
38 thought that using the space to analyze the  
39 action selection process in Animat would  
40 be more relevant.



41 **Reviewer 3:**

42 - We will make sure to clarify on the intuition behind our objective. Assuming an agent improves the expected return of the policy  
43 with each episode, maximizing the cumulative return is equivalent to maximizing the AUC. Intuitively, when plotted, the quicker an  
44 agent reaches an optimal policy, the larger the AUC will be. Because the only variable that is optimized is the exploration policy, it  
45 learns to exploit structures present on the tasks to achieve this.

46 - The behaviors learned for pole balancing and self-driving task were a bit more intuitive. In pole balancing, the exploration policy  
47 learns that if the pole is about to fall to the right, the cart should compensate by moving to the right, and vice versa. Similarly, in  
48 the self-driving task, if the car is driving to much to the right the appropriate action should be moving to the left by some amount  
49 and vice versa. In general, certain actions can safely be omitted from exploration in certain states. We will make a more detailed  
50 discussion regarding this result.

51 - Figure 5 is showing that, in the first half of the set of tasks (while the advisor has not trained sufficiently), the exploration policy is  
52 inefficient. The second half of set of tasks demonstrate that the exploration policy is improving over the initial exploration policy and  
53 leads to improving learning over random exploration. We will clarify the points discussed in section 6.1

54 - Thank you for pointing out REPTILE as an alternative meta learning method.