

1 We thank the reviewers for the thorough feedbacks. Based on those, we have made numerous improvements.

2 **Implement a new IM baseline: ICM (Pathak 2017 [23].**  
 3 **Original code is for discrete actions.)** As suggested by re-  
 4 viewer #1, #3, we have implemented ICM for the synthetic  
 5 environment (Sec.4, Fig. 3 of the manuscript). The ICM base-  
 6 line uses SAC with an augmented reward:  $r_t = r_t^{\text{ex}} + \alpha r_t^{\text{in}}$ ,  
 7 where  $r_t^{\text{ex}}$  is the extrinsic reward (negative distance to goal) and  
 8  $r_t^{\text{in}}$  is an intrinsic reward.

9 The first experiment (Fig. 1 Left) follows the original ICM,  
 10 where the intrinsic reward signal is given by the total predic-  
 11 tion error:  $r^{\text{in}} = \sum_i e_i(t)$ , where the sum is over all goal  
 12 spaces/coordinates. Furthermore, we adapted ICM to make use  
 13 of the surprise signals that have shown to be important in the  
 14 manuscript. Thus, in a second experiment (Fig. 1 Right), the  
 15 intrinsic reward is given by the surprise signal:  $r^{\text{in}} = \max_i \text{surprise}_i(t)$ , where max is over goal  
 16 spaces. Despite scanning the hyperparameter  $\alpha$ , both IM baselines perform poorly and only solve the locomotion task, see Fig. 1.  
 17 Despite the seemingly simple environment, a random encounter of objects in continuous control is rare, given an agent  
 18 with heavy mass and a large arena.

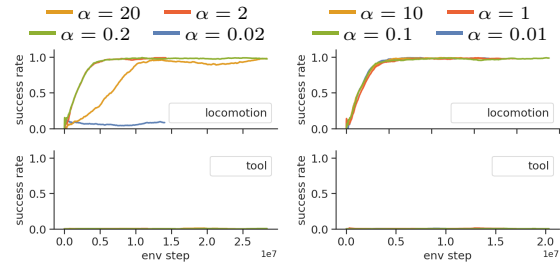


Fig. 1: Synthetic environment in Sec. 4. Left: prediction error; right: surprise.  $\alpha$  is a hyperparameter we scanned for.

19 To address rev. #2’s concern over “object can’t be moved, a model-error driven IM will stop”, we first clarify that the  
 20 issue, in fact, lies with the “random object” (in Sec. 4), not an unmovable object. We further tested the above-mentioned  
 21 IM baseline with the random object. The plot is similar to “tool” in Fig. 1 and we omit it due to space constraints.

22 **Clarify novelty and main contributions** We agree that each individual component is not original, as we have clearly  
 23 indicated they are from task-motion planning, IM, RL communities. We have already given references in the manuscript  
 24 (including Klyubin and Battaglia’s work(s) mentioned by rev. #1). But combining them to successfully solve the  
 25 continuous control and robot trajectory optimization problem is novel (cf. rev. #3, originality).

26 Rev. #1 suggested that the environments could be solved by classic planning methods. If one has an environment model  
 27 with an analytically (or accurate numerical) gradient, iLQR(G) may (without guarantee) solve the nonlinear program  
 28 (NLP). We have discussed this and other planning ideas (e.g. PRM) in the related work section. However, this paper is  
 29 based on model-free RL to solve the robot trajectory optimization through contact. We demonstrated IM/RL can solve  
 30 this as an alternative to NLP/sampling-based planning. This is beyond the scope of existing works such as Klyubin et al.

31 It is true that our method shares certain points with the concept of empowerment. We would like to emphasize that the  
 32 structure that we proposed leads to more efficient learning while maintaining the idea of maximizing controllability.

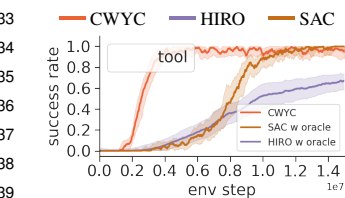


Fig. 2: Baselines with reward-shaping

33 **Concerning the complexity of our method** We acknowledge that the original Fig. 1  
 34 suggests an overwhelming complexity due to the detailed break-down (we will simplify  
 35 this). In fact, our *inductive bias* (c.f. [Tenenbaum (2011) “How to grow...”]) has only  
 36 3 modules (not 8): the task selector, planner, and subgoal generator. All other modules  
 37 are common among RL algorithms. In the ablation studies, we demonstrated that every  
 38 component is required to solve the task/maintain data efficiency. To further validate  
 39 this claim, we report additional results in Fig. 2, where the baselines are able to learn  
 40 the tool task with a hand-engineered reward:  $r_t = r_t^{\text{ex}} - \text{dist}(\text{agent-pos}_t, \text{tool-pos}_t)$ .  
 41 Therefore, our method in fact removes this additional layer of supervision.

42 **Further improvements.** Code is uploaded to the website as given in the paper. Con-  
 43 cerning our argument for playfulness, see [Smith (2005) “The dev. of embodied...”; Ryan (2000) “Intrinsic and  
 44 extrinsic...”]. Regarding prediction error vs. learning progress: prediction error fails in stochastic environments, see  
 45 [Oudeyer (2007) “Intrinsic motiv. systems...”; Burda (2018) “Large-scale...”].

46 **Q:** subgoal attention requires attending over all possible goals...? **A:** Our specific form of the goal generation network  
 47 allows for a closed-form solution to compute the argmax of the function. **Q:** The task graph is not a function of the  
 48 particular goal in the final task...? **A:** True. A limitation of our current architecture. **Q:** Goals within one task have  
 49 different difficulty. **A:** True. Interesting future direction. **Q:** When is the transition between sub-tasks happening...? **A:**  
 50 Your understanding is correct. If the goal can not be reached, the rollout is terminated after the maximum timesteps per  
 51 rollout is reached. We clarify this.

52 All text errors or vague language will be fixed. We have addressed other review comments but omit reporting them here  
 53 due to the space constraint. We gratefully acknowledge your help in improving the work.