1 We thank the reviewers for clear and thoughtful feedback, and respond to specific points raised by reviewers below.

2 **R2:** "how the approach compares to [22]." **R3:** "absence of prior work that it out-performs".

3 To address the primary concerns of **R2** and **R3**, we present results of new comparisons to Gupta et al. [22] on the Fixed
4 ViZDoom experimental setting in Table 1. This comparison ([22]) is representative of "train[ing] an agent and task
5 distribution using one of the 10s of DIAYN-like approaches" (**R3**) before freezing the task distribution and running
6 meta-learning in a "pipelined" manner. However, we note that [22] considers environments with simpler, ground-truth
7 state, as opposed to pixel observations.

Table 1: Comparing to [22].

|  | Avg. Succ. |
|---|---|
| **i.** [22] | 0.291 |
| **ii. Ours**, pipeline | 0.535 |
| **iii.** [22], smart-init | 0.405 |
| **iv. Ours**, full | **0.625** |

8 The compared approaches are: **(i)** [22], which uses DIAYN [13] for task acquisition,
9 adapted for pixel observations; **(ii)** an ablation of our method – "pipelined CARML"
10 – more similar to [22], for an apples-to-apples comparison; **(iii)** [22], but initializing
11 the DIAYN discriminator of with the image encoder of **(ii)**, to address failure modes
12 of applying [22] in visual domains; and **(iv)** CARML, our full method.

13 **Our approach outperforms [22] on transfer to test tasks.** The benefit of our task acquisition method over that of
14 DIAYN (which [22] uses) is indicated by the improvement from **(i)** and **(iii)** to **(ii)**. The benefit of using a curriculum
15 for meta-learning over the pipelined approach of [22] is indicated by the improvement from **(ii)** to **(iv)**. **Please find**
16 **discussion of these results at the end of the page.** We will include these and further experiments on the remaining
17 settings in our revision.

18 **R3**: "Show that ... the newly proposed task is super useful".

19 We note that the environments considered are nearly identical to the navigation setting of [55] (though ours is more
20 challenging insofar as no task description is given) and the manipulation setting used in [34], among others. Our
21 work is among the first to study unsupervised meta-RL in visual domains, addressing challenges of pixel observation
22 trajectories and partial observability, among others, which exacerbate the challenges of unsupervised RL and meta-RL.

23 **Populating** $\mathcal{D}$ **(R2).** We choose the simplest strategy that keeps complexity constant: sample a fixed number of
24 trajectories uniformly at random from the entire history, i.e. reservoir sampling. We used a reservoir of 1000 trajectories
25 (not tuned). We agree with **R2** that more sophisticated sampling strategies are worth pursuing in future work.

26 **Comparison Details.** Differing from [22], we use $RL^2$ instead of MAML for more direct comparability; to our
27 knowledge, policy gradient MAML has yet to be successfully implemented in RL domains with pixel observations.
28 Comparison **(ii)** uses a contextual policy to co-adapt with the task distribution before freezing the task distribution
29 and meta-learning with $RL^2$. Results are reported for transfer to the Fixed ViZDoom test tasks, analogous to results
30 in Figure 5a of submission. We use the same hyper-parameters for skill acquisition (i.e. number of skills) as existing
31 experiments. In Table 1, we report the average of two runs per approach, but will use more in our revision.

32 **Comparison Discussion (R2, R3).** We find the task acquisition of DIAYN variants **(i, iii)** to suffer from an effect akin
33 to mode-collapse; the policy's data distribution collapses to a smaller subset of the trajectory space (one or two modes),
34 and tasks correspond to minor variations of these modes. Skill acquisition methods such as DIAYN rely purely on
35 discriminability of states/trajectories under skills, which can be more easily satisfied in high-dimensional observation
36 spaces and can thus lead to such mode-collapse (related to the instability of GAN methods noted by **R1**). Moreover,
37 they do not a provide a direct mechanism for furthering exploration once skills are discriminable.

38 On the other hand, the proposed task acquisition approach (Alg. 2, Sections 3.2, 3.4) **fits a generative model over**
39 **jointly learned discriminative features**, and is thus not only **less susceptible to mode-collapse** (w.r.t the policy data
40 distribution), but also allows for density-based exploration (Section 3.3). Indeed, we find that **(iii)** seems to mitigate
41 mode-collapse – benefiting from a pretrained encoder from **(ii)** – but does not entirely prevent it. Overall, in terms of
42 meta-transfer to hand-crafted test tasks, the DIAYN variants **(i, iii)** perform worse than pipelined CARML **(ii)**, due to
43 the poorer diversity in the task distribution. We will incorporate this comparison, as well as additional visualizations
44 (i.e. skill maps) of all discussed methods, in the revised Appendix.

45 Moreover, **(ii)** performs worse than "full CARML" **(iv)**. As in the paper, we hypothesize that this is due to the challenge
46 of meta-learning more complex task distributions – compared to full CARML, the distribution of trajectories eventually
47 discovered by the contextual policy of **(ii)** may be just as diverse and structured, but meta-learning the corresponding
48 task distribution directly from scratch is harder. This shows **the benefit of co-adapting tasks with the meta-learner**
49 **(iv)** as opposed to using a separate agent **(ii)**, and the value of investigating the effects of curricula on meta-learning.

50 **[13]** Eysenbach, Gupta, Ibarz, Levine. Diversity is all you need: learning skills without a reward function. ICLR 2019.
51 **[22]** Gupta, Eysenbach, Finn, Levine. Unsupervised meta-learning for reinforcement learning. arXiv:1806.04640, 2018.
52 **[34]** Nair, Pong, Dalal, Bahl, Lin, and Levine. Visual reinforcement learning with imagined goals. NeurIPS 2018.
53 **[55]** Chaplot, Sathyendra, Pasumarthi, Rajagopal, Salakhutdinov. Gated-attention architectures for task-oriented
54 language grounding. AAAI 2018.