

# UNSUPERVISED META-LEARNING FOR REINFORCEMENT LEARNING

**Anonymous authors**

Paper under double-blind review

## ABSTRACT

Meta-learning algorithms learn to acquire new tasks more quickly from past experience. In the context of reinforcement learning, meta-learning algorithms can acquire reinforcement learning procedures to solve new problems more efficiently by utilizing experience from prior tasks. The performance of meta-learning algorithms depends on the tasks available for meta-training: in the same way that supervised learning generalizes best to test points drawn from the same distribution as the training points, meta-learning methods generalize best to tasks from the same distribution as the meta-training tasks. In effect, meta-reinforcement learning offloads the design burden from algorithm design to task design. If we can automate the process of task design as well, we can devise a meta-learning algorithm that is truly automated. In this work, we take a step in this direction, proposing a family of unsupervised meta-learning algorithms for reinforcement learning. We motivate and describe a general recipe for unsupervised meta-reinforcement learning, and present an instantiation of this approach. Our conceptual and theoretical contributions consist of formulating the unsupervised meta-reinforcement learning problem and describing how task proposals based on mutual information can in principle be used to train optimal meta-learners. Our experimental results indicate that unsupervised meta-reinforcement learning effectively acquires accelerated reinforcement learning procedures without the need for manual task design and significantly exceeds the performance of learning from scratch.

## 1 INTRODUCTION

Reusing past experience for faster learning of new tasks is a key challenge for machine learning. Meta-learning methods propose to achieve this by using past experience to explicitly optimize for rapid adaptation (Mishra et al., 2017; Snell et al., 2017; Schmidhuber, 1987; Finn et al., 2017a; Duan et al., 2016b; Gupta et al., 2018; Wang et al., 2016; Al-Shedivat et al., 2017; Rakelly et al., 2019). In the context of reinforcement learning, meta-reinforcement learning algorithms can learn to solve new reinforcement learning tasks more quickly through experience on past tasks (Duan et al., 2016b; Gupta et al., 2018). Typical meta-reinforcement learning algorithms assume the ability to sample from a pre-specified task distribution, and these algorithms learn to solve new tasks *drawn from this distribution* very quickly. However, specifying a task distribution is tedious and requires a significant amount of supervision (Finn et al., 2017b; Duan et al., 2016b) that may be difficult to provide for large real-world problem settings. The performance of meta-learning algorithms critically depends on the meta-training task distribution, and meta-learning algorithms generalize best to new tasks which are drawn from the same distribution as the meta-training tasks (Finn & Levine, 2018). In effect, meta-reinforcement learning offloads some of the design burden from algorithm design to designing a sufficiently broad and relevant distribution of meta-training tasks. While this greatly helps in acquiring representations for fast adaptation to the specified task distribution, a natural question is whether we can do away with the need for manually designing a large family of tasks, and develop meta-reinforcement learning algorithms that learn only from unsupervised environment interaction. In this paper, we take an initial step toward the formalization and design of such methods.

Our goal is to automate the meta-training process by removing the need for hand-designed meta-training tasks. To that end, we introduce unsupervised meta-reinforcement learning: meta-learning from a task distribution that is acquired automatically, rather than requiring manual design of the meta-training tasks. Unsupervised meta-reinforcement learning methods must solve two difficult

problems together: meta-reinforcement learning with broad task distributions, and unsupervised exploration for proposing a wide variety of tasks for meta-learning. Since the assumptions of our method differ fundamentally from prior meta-reinforcement learning methods (we do not assume access to hand-specified meta-training tasks), the best points of comparison for our approach are learning meta-test tasks entirely from scratch with conventional reinforcement learning algorithms. Our method can also be thought of as a data-driven *environment-specific* initialization procedure for deep neural network policies, somewhat related to data-driven initialization procedures explored in supervised learning (Krähenbühl et al., 2015). However, as indicated by Finn & Levine (2017), this procedure goes beyond simply being an initialization, and essentially learns an entire learning algorithm that is as expressive as any recurrent meta-learner.

The primary contributions of our work are to propose a framework for unsupervised meta-reinforcement learning, sketch out a family of unsupervised meta-reinforcement learning algorithms, provide a theoretical derivation that allows us to reason about the optimality of unsupervised meta-reinforcement learning methods in terms of mutual information, and describe an instantiation of an algorithm from this family that builds on a recently proposed procedure for unsupervised exploration (Eysenbach et al., 2018) and model-agnostic meta-learning (MAML) (Finn et al., 2017a). In addition to our theoretical derivations, we provide an empirical evaluation that studies the performance of two variants of our approach on simulated control tasks. Our experimental evaluation shows that, for a variety of tasks, unsupervised meta-reinforcement learning can effectively acquire reinforcement learning procedures that perform significantly better than standard reinforcement learning methods that learn from scratch, without additional task knowledge.

## 2 RELATED WORK

Our work lies at the intersection of meta reinforcement learning, goal generation, and unsupervised exploration. Meta-learning algorithms use data from multiple tasks to learn how to learn, acquiring rapid adaptation procedures from experience (Schmidhuber, 1987; Naik & Mammone, 1992; Thrun & Pratt, 1998; Bengio et al., 1992; Hochreiter et al., 2001; Santoro et al., 2016; Andrychowicz et al., 2016; Li & Malik, 2017; Ravi & Larochelle, 2017; Finn et al., 2017a; Munkhdalai & Yu, 2017; Snell et al., 2017). These approaches have been extended into the setting of reinforcement learning (Duan et al., 2016b; Wang et al., 2016; Finn et al., 2017a; Sung et al., 2017; Mishra et al., 2017; Gupta et al., 2018; Mendonca et al., 2019; Houthoofd et al., 2018; Stadie et al., 2018), though their performance in practice depends on the user-specified meta-training task distribution. We aim to lift this limitation, and provide a general recipe for avoiding manual task engineering for meta-RL. To that end, we make use of unsupervised task proposals. These proposals can be obtained in a variety of ways, including adversarial goal generation (Sukhbaatar et al., 2017; Held et al., 2017), information-theoretic methods (Gregor et al., 2016; Eysenbach et al., 2018; Co-Reyes et al., 2018; Achiam et al., 2018), and even random functions. We argue that, theoretically, methods based on mutual information have the potential to provide *optimal* task proposals for unsupervised meta-reinforcement learning. Exploration methods that seek out novel states are also closely related to goal generation methods (Pathak et al., 2017; Schmidhuber, 2009; Bellemare et al., 2016; Osband et al., 2016; Stadie et al., 2015), but do not by themselves aim to generate new tasks or learn to adapt more quickly to new tasks, only to achieve wide coverage of the state space. These methods are complementary to our approach, but address a distinct problem.

Related to our work are prior methods that study the training of goal-conditioned policies (Schaul et al., 2015; Pong et al., 2018; Andrychowicz et al., 2017). Indeed, our theoretical derivation first studies the goal reaching case, before generalizing it to the more general case of arbitrary tasks. This generalization allows unsupervised meta-learning methods to solve arbitrary tasks at meta-test time without being restricted to a task parameterization. Additionally, we discuss why the framework of meta-learning gives us theoretical benefits over the goal reaching paradigm.

## 3 UNSUPERVISED META-REINFORCEMENT LEARNING

The goal of unsupervised meta-reinforcement learning is to take an environment and produce a learning algorithm specifically tailored to *this* environment that can quickly learn to maximize reward on *any* task reward in this environment. This learning algorithm should be meta-learned without

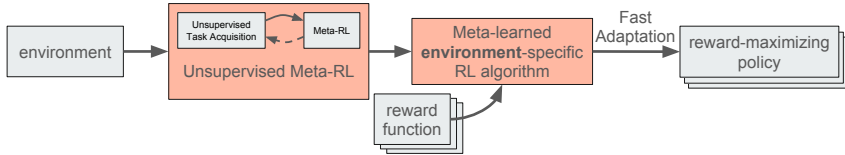


Figure 1: **Unsupervised meta-reinforcement learning:** Given an environment, unsupervised meta-reinforcement learning produces an environment-specific learning algorithm that quickly acquire new policies that maximizes any task reward function.

requiring *any* human supervision. We can formally define unsupervised meta-reinforcement learning in the context of a controlled Markov process (CMP) – a Markov decision process without a reward function,  $C = (S, A, P, \gamma, \rho)$ , with state space  $S$ , action space  $A$ , transition dynamics  $P$ , discount factor  $\gamma$  and initial state distribution  $\rho$ . Our goal is to learn a learning algorithm  $f$  on this CMP, which can subsequently learn new tasks efficiently in this CMP for a new reward function  $R_i$ . The CMP along with this reward function  $R_i$  produces a Markov decision processes  $M_i = (S, A, P, \gamma, \rho, R_i)$ . The goal of the learning algorithm  $f$  is to learn an optimal policy  $\pi_i^*(a|s)$  for *any* reward function  $R_i$  that is provided with the CMP. Crucially,  $f$  must be learned without access to any reward functions  $R_i$ , using only unsupervised interaction with the CMP. The reward is only provided at meta-test time. The implicit assumption in this formulation is that different tasks at test-time will all be using the same dynamics but with different reward functions. In this section, we will first sketch out a general recipe for an unsupervised meta-reinforcement learning algorithm, then present a derivation for an *optimal* unsupervised meta-learning method, and then instantiate a practical approximation to this theoretically motivated approach using components from recently proposed exploration and meta-learning algorithms.

### 3.1 A GENERAL RECIPE

Our framework unsupervised meta-reinforcement learning consists of a task proposal mechanism and a meta-learning method. Formally, we will define the task distribution as a mapping from a latent variable  $z \sim p(z)$  to a reward function  $r_z(s, a) : \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}$ . That is, for each value of the random variable  $z$ , we have a different reward function  $r_z(s, a)$ . Under this formulation, learning a task distribution amounts to optimizing a parametric form for the reward function  $r_z(s, a)$  that maps each  $z \sim p(z)$  to a different reward function. The choice of this parametric form represents the most important decision for an unsupervised meta-learning method, and we will discuss a theoretical framework that allows us to make this choice in the following section so as to minimize worst case regret of the subsequently meta-learned learning algorithm  $f$ . The second is the meta-learning method, which takes the family of reward functions induced by  $p(z)$  and  $r_z(s, a)$ , and meta-learns a reinforcement learning algorithm  $f$  that can quickly adapt to any task from the task distribution defined by  $p(z)$  and  $r_z(s, a)$ . The meta-learned algorithm  $f$  can learn new tasks quickly at meta-test time, when a user-specified reward function is actually provided. This generic design for an unsupervised meta-reinforcement learning algorithm is summarized in Figure 1.

The “no free lunch theorem” (Wolpert et al., 1995; Whitley & Watson, 2005) might lead us to expect that a truly generic approach to learning a task distribution would not yield a learning procedure  $f$  that is effective on any real tasks – or even on the meta-training tasks. Note, however, that an unsupervised meta-learning algorithm *can* in fact collect and organize meaningful information about the environment (that is, about the CMP) even without a reward function. Therefore, the capability of unsupervised meta-learning to improve over a learner that learns each new task from scratch depends on the degree to which it can acquire useful knowledge about the task. In the following sections, we will discuss how we can formulate an optimal unsupervised meta-learner that minimizes regret on new meta-test tasks in the absence of any prior knowledge, and then show how we can instantiate an approximation to this theoretically motivated method.

### 3.2 OPTIMAL META-LEARNERS

In order to perform our derivation, we first need to define an abstract notion of an optimal meta-learner, given a task distribution. We assume that an optimal meta-learner takes in a distribution over tasks and outputs a learning procedure  $f$  that minimizes expected regret when learning tasks drawn from

the *same* distribution as meta-training. As before, the task distribution is defined by a latent variable  $z \sim p(z)$  and a reward function  $r_z(s, a) : \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}$ . The objective that a optimal meta-learner optimizes is

$$\min_f \mathbb{E}_{z \sim p(z)} [\text{REGRET}(f, z, r_z(s, a))], \quad (1)$$

where the regret is measured during adaptation to a new task corresponding to  $z$ . This is equivalent to the expected reward objective used by most meta-reinforcement learning methods Finn et al. (2017a); Duan et al. (2016b). As we show below, the optimal meta-learner has distinct behavior conditional on the specific task distribution, and can be characterized for particular classes of task distributions. Given this definition of an optimal meta-learner, we consider how we can construct unsupervised meta-learning algorithms.

### 3.3 SPECIAL CASE: GOAL-REACHING TASKS

We will first derive an optimal unsupervised meta-learner for the special case of goal reaching tasks, and then generalize this to the case of all possible tasks in Section 3.4. Specifically, we consider episodes with finite horizon  $T$ , and discount factor of  $\gamma = 1$ . Tasks correspond to reaching an unknown goal state  $s_g$ . We will only consider the agent’s state at the last time step in each episode, so the (unknown) reward function is always of the form

$$r_g(s_t) \triangleq \mathbb{1}(t = T) \cdot \mathbb{1}(s_t = g).$$

We will first assume that goal states are drawn from some distribution  $p(s_g)$ , and later will show how we can remove this assumption. We define  $\rho_\pi^T(s)$  as the probability that policy  $\pi$  visits state  $s$  at time step  $t = T$ . If  $s_g$  is the true goal, then the event that the policy  $\pi$  reaches  $s_g$  at the final step of an episode is a Bernoulli random variable with parameter  $p = \rho_\pi^T(s_g)$ . Thus, the expected *hitting time* of this goal state is

$$\text{HITTINGTIME}_\pi(s_g) = \frac{1}{\rho_\pi^T(s_g)}$$

Recall that the goal state  $s_g$  is unknown. The cumulative regret is therefore the expectation of the hitting time, taken with respect to  $p(s_g)$ :

$$\text{REGRET}_p(\pi) = \int \text{HITTINGTIME}_\pi(s_g) p(s_g) ds_g = \int \frac{p(s_g)}{\rho_\pi^T(s_g)} ds_g \quad (2)$$

In this special case, an optimal meta-learner as defined in Section 3.2 will explore for a number of episodes until it finds the goal state. After the meta-learner finds the goal state, it would always choose the trajectory that reaches that goal state under deterministic dynamics. Thus, the cumulative regret of the meta-learner is the number of episodes required to find the goal state. By our assumption that the meta-learner could only learn if it had reached the goal state at the end of the episode, the meta-learner cannot use learn information about multiple goals within a single episode. We can minimize the regret in Equation 2 w.r.t. the marginal distribution  $\rho_\pi^T$  using the calculus of variations, which tells us that the (exploration) policy for the optimal meta-learner,  $\pi^*$ , satisfies:

$$\rho_{\pi^*}^T(s_g) = \frac{\sqrt{p(s_g)}}{\int \sqrt{p(s'_g)} ds'_g} \quad (3)$$

The analysis so far tells us how to obtain the optimal meta-learner if we were given the goal sampling distribution,  $p(s_g)$ . If we do not know this distribution, then we cannot compute the optimal policy using Equation 3. In this case, we resort to bounding the *worst-case* regret of our policy:

$$\min_\pi \max_p \text{REGRET}_p(\pi) \quad (4)$$

**Lemma 1.** *Let  $\pi$  be a policy for which  $\rho_\pi^T(s)$  is uniform. Then  $\pi$  has lowest worst-case regret.*

The proof is in Appendix B. Given this result, we know that the optimal meta-learner should have a uniform state marginal distribution. The minimax optimal meta-learner corresponds to a uniform distribution over goal states, so we can acquire this meta-learner by training on a goal-reaching task distribution where the goals are uniformly distributed. Manually crafting this distribution is hard, especially in high-dimensional settings. We are therefore left with the problem of devising a

method that can propose goals to our optimal meta-learner during unsupervised meta-training that are distributed uniformly. Recall that the task proposal mechanism is defined in terms of a latent variable  $z \sim p(z)$  and a parameterized reward function  $r_z(s, a)$ . For each  $z \sim p(z)$ , let us also introduce the optimal policy for the corresponding reward, which we will denote  $\pi(a | s, z)$ . This policy is *not* the meta-learned model, but simply a component of the task proposal mechanism. This policy induces a distribution over terminal states,  $p(s_T | z)$ . We will show that a uniform goal proposal distribution is obtained by a goal proposal mechanism that maximizes the mutual information between  $z$  and the final state  $s_T$ :

$$I(s_T; z) \triangleq \mathcal{H}[s_T] - \mathcal{H}[s_T | z] \quad (5)$$

Observe that this objective contains two competing terms. The first term,  $\mathcal{H}[s_T]$ , says that  $\pi(a | s, z)$ , when aggregated over many latents  $z$ , has a high-entropy state distribution. This term is maximized when the state distribution is uniform, as desired. The second term,  $\mathcal{H}[s_T | z]$ , says that  $\pi(a | s, z)$ , when conditioned on a particular latent  $z$ , should go to only a few states. If we optimize this objective, we obtain a marginal distribution over final states that is uniform (proof in Appendix B):

**Lemma 2.** *Mutual information  $I(s_T; z)$  is maximized by a task distribution  $p(s_g)$  that is uniform over goal states.*

We can then recover a final time step reward function for each value of  $z$  as  $r_z(s_T, a_T) \triangleq \log p(s_T | z)$ , where the probability distribution is taken with respect to the optimal policy  $\pi(a | s, z)$ . One peculiar fact of the goal reaching tasks is that agents cannot take epistemic (i.e., information gathering) actions, since the agent only acquires information about the goal state when it reaches it. In this particular setting, meta-learners that perform posterior sampling for exploration, like the one described above, are optimal. However, in more general settings, meta-learning algorithms can also take epistemic actions, potentially performing better than posterior sampling.

### 3.4 GENERAL CASE: TRAJECTORY-MATCHING TASKS

To extend the analysis in the previous section to the general case, and thereby derive a framework for optimal unsupervised meta-learning, we will consider “trajectory-matching” tasks. These tasks are a trajectory-based generalization of goal reaching: while goal reaching tasks only provide a positive reward when the policy reaches the goal state, trajectory-matching tasks only provide a positive reward when the policy executes the optimal trajectory.<sup>1</sup> These non-Markovian tasks essentially amount to a problem where an RL algorithm must “guess” the optimal policy, and only receives a reward if its behavior is perfectly consistent with that optimal policy. We will show that mutual information between  $z$  and *trajectories* yields the minimum regret solution in this case, and then show that unsupervised meta-learning for the trajectory-matching task is at least as hard as unsupervised meta-learning for general tasks (though in practice it is considerably harder).

Formally, we define a distribution of trajectory-matching tasks by a distribution over goal trajectories,  $p(\tau^*)$ . For each goal trajectory  $\tau^*$ , the corresponding trajectory-level reward function is

$$r_{\tau^*}^*(\tau) \triangleq \mathbb{1}(\tau = \tau^*) \quad (6)$$

As before, we define the hitting time as the expected number of episodes to match the target trajectory:

$$\text{HITTINGTIME}_{\pi}(\tau^*) = \frac{1}{\pi(\tau^*)} \quad (7)$$

We then define regret as the expected hitting time:

$$\text{REGRET}_p(\pi) = \int \text{HITTINGTIME}_{\pi}(\tau) p(\tau) d\tau = \int \frac{p(\tau)}{\pi(\tau)} d\tau \quad (8)$$

Using the same derivation as before, the exploration policy for the optimal meta-learner is

$$\pi^*(\tau) = \frac{\sqrt{p(\tau)}}{\int \sqrt{p(\tau')} d\tau'}. \quad (9)$$

However, obtaining such a policy requires knowing the trajectory distribution  $p(\tau)$ . In the setting where  $p(\tau)$  is unknown, the minimax policy is simply uniform:

<sup>1</sup>In a stochastic CMP, this corresponds to all trajectories that can be executed with non-zero probability by the deterministic optimal policy for the unknown reward function, but we will present the derivation for the deterministic case for simplicity.

**Lemma 3.** *Let  $\pi$  be a policy for which  $\pi(\tau)$  is uniform. Then  $\pi$  has lowest worst-case regret.*

How can we acquire a policy with a uniform trajectory distribution? Repeating the steps above, we learn a collection of skills using a *trajectory-level* mutual information objective:

$$I(\tau; z) = \mathcal{H}[\tau] - \mathcal{H}[\tau | z] \quad (10)$$

Using the same reasoning as Section 3.3, the optimal policy for this objective has a uniform marginal distribution over trajectories that, conditioned on a particular latent  $z$ , deterministically produces a single trajectory in a deterministic CMP, or one out of a set of trajectories in a stochastic CMP. As in Section 3.3, we define a distribution over reward functions as  $r_z(\tau) \triangleq \log p(\tau | z)$ . At optimality, each  $z$  corresponds to exactly one trajectory  $\tau_z$ , so the reward function  $r_z(\tau)$  simply indicates whether  $\tau$  is equal to  $\tau_z$ . Recalling that the marginal distribution over trajectories  $\int p(\tau | z)p(z)dz$  is uniform at optimum, the distribution of reward functions  $r_z$  corresponds to a uniform distribution over trajectories. Thus, meta-learning on the rewards from trajectory-level mutual information results in the minimax-optimal meta-learner.

It is important to note here that the specific case of trajectory-matching is a super-set of the problem of optimizing any possible Markovian reward function at test-time. For a given initial state distribution, each reward function is optimized by a particular trajectory. However, trajectories produced by a non-Markovian policy (i.e., a policy with memory) are not necessarily the unique optimum for any Markovian reward function. Let  $R_\tau$  denote the set of trajectory-level reward functions, and  $R_{s,a}$  denote the set of all state-action level reward functions. Bounding the worst-case regret on  $R_\tau$  minimizes and upper bound on the worst-case regret on  $R_{s,a}$ :

$$\min_{r_\tau \in R_\tau} \mathbb{E}_\pi [r_\tau(\tau)] \leq \min_{r \in R_{s,a}} \mathbb{E}_\pi \left[ \sum_t r(s_t, a_t) \right] \quad \text{for all policies } \pi.$$

In general, this bound is loose because the set of all Markovian reward functions is smaller than the set of all trajectories. However, this bound becomes tight when considering meta-learning on the set of all possible (non-Markovian) reward functions.

In the discussion of meta-learning thus far, we have considered tasks where the reward is provided at the last time step  $T$  of each episode. In this particular case, the best that an optimal meta-learner can do is go directly to a goal or execute a particular trajectory at every episode according to the optimal exploration policy:  $\rho_{\pi^*}^T(s_g) = \frac{\sqrt{p(s_g)}}{\int \sqrt{p(s'_g)} ds'_g}$  for goal reaching or  $\pi^*(\tau) = \frac{\sqrt{p(\tau)}}{\int \sqrt{p(\tau')} d\tau'}$  for trajectory matching. All intermediate states in the trajectory are uninformative, thus making instances of meta-learning algorithms which explore via schemes like posterior sampling optimal for this class of problems. In the more general case with arbitrary reward functions, intermediate rewards along a trajectory may be informative, and the optimal exploration strategy may be somewhat different (Rothfuss et al., 2019; Duan et al., 2016b; Wang et al., 2016). However, the analysis presented in Section 3.4 provides us insight into the behavior of optimal meta-learning algorithms and allows us to understand the qualities desirable for unsupervised task proposals.

Through our analysis, we introduce the notion of optimal meta-learners and analyze their exploration behavior and regret on a class of goal reaching problems. We show that on these problems, when the test-time task distribution is unknown, the optimal meta-training task distribution for minimizing worst-case test-time regret is *uniform* over the space of goals. We show that this optimal task distribution can be acquired by a simple mutual information maximization scheme. We subsequently extend the analysis to the more general case of matching arbitrary trajectories, as a proxy for the more general class of arbitrary reward functions. In the following section we discuss how we can derive a practical algorithm for unsupervised meta-learning from this analysis, and empirically validate it’s effectiveness.

### 3.5 A PRACTICAL ALGORITHM

Following the derivation in the previous section, we can instantiate a practical unsupervised meta-reinforcement learning algorithm by constructing a task proposal mechanism based on a mutual information objective. To this end, we adopt the DIAYN algorithm (Eysenbach et al., 2018). DIAYN constructs tasks by approximately maximizing a mutual information objective. A variety of different

mutual information objectives can be formulated, including mutual information between single states and  $z$  (Eysenbach et al., 2018), pairs of start and end states and  $z$  Gregor et al. (2016), and entire trajectories and  $z$  Achiam et al. (2018). While the latter objective is motivated by our theory, we found the simple state-based objective to work well in practice, and leave a full examination of possible mutual information objectives for future work. DIAYN optimizes mutual information by training a *discriminator* network  $D_\phi(z|\cdot)$  that predicts which  $z$  was used to generate the states in a given rollout according to a *latent-conditioned policy*  $\pi(a|s, z)$ . Depending on what the discriminator is conditioned on, we obtain different mutual information objectives. We use  $D_\phi(z|s)$ , resulting in the objective  $I(s; z)$  (Eysenbach et al., 2018). The reward in DIAYN is given by  $r_z(s, a) = \log(D_\phi(z|s))$ .

The complete unsupervised meta-learning algorithm follows the recipe in Figure 1: first, we acquire  $r_z(s, a)$  by running DIAYN, which learns  $D_\phi(z|s)$  and a latent-conditioned policy  $\pi(a|s, z)$  (which is discarded). Then, we use  $z \sim p(z)$  to propose tasks  $r_z(s, a)$  to a standard meta-reinforcement learning algorithm. We use MAML Finn et al. (2017a), though any other meta-reinforcement learning method could be used in principle. This method is summarized in Algorithm 1.

---

**Algorithm 1:** Unsupervised Meta-Reinforcement Learning Pseudocode

---

**Data:**  $\mathcal{M} \setminus R$ , an MDP without a reward function

**Result:** a learning algorithm  $f : D \rightarrow \pi$

Initialize  $D = \emptyset$

$D_\phi \leftarrow \text{DIAYN}()$  or  $D_\phi \leftarrow \text{random}$

**while not converged do**

    Sample latent task variables  $z \sim p(z)$

    Extract corresponding task reward functions  $r_z(s)$  using  $D_\phi(z|s)$

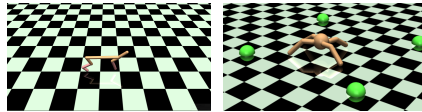
    Update  $f$  using MAML with reward  $r_z(s)$

---

In addition to mutual information maximizing task proposals, we will also consider random task proposals, where we also use a discriminator as the reward, according to  $r(s, z) = \log D_{\phi_{rand}}(z|s)$ , but where the parameters  $\phi_{rand}$  are chosen randomly (i.e., a random weight initialization for a neural network). While such random reward functions are not optimal, we find that they can be used to acquire useful task distributions for simple tasks, though they are not as effective as the tasks become more complicated. This empirically reinforces the claim that unsupervised meta-reinforcement learning does not in fact violate any “no free lunch” principle – even simple task proposals that cause the meta-learner to explore the CMP can already accelerate learning of new tasks.

## 4 EXPERIMENTAL EVALUATION

In our experiments, we aim to understand whether unsupervised meta-learning as described in Section 3.1 can provide us with an accelerated reinforcement learning procedure on new tasks. Whereas standard meta-learning requires a hand-specified task distribution at meta-training time, unsupervised meta-learning learns the task distribution through unsupervised interaction with the environment. A fair baseline that likewise uses *no reward supervision* at training time, and only uses reward for learning new tasks is learning via RL from scratch without any meta-learning. As an upper bound, we include the *unfair* comparison to a standard meta-learning approach, where the meta-training distribution is manually designed. This method has access to a hand-specified task distribution that is not available to our method. We evaluate two variants of our approach: (a) task acquisition based on DIAYN followed by meta-learning using MAML, and (b) task acquisition using a randomly initialized discriminator followed by meta-learning using MAML.



Half-Cheetah

Ant

### 4.1 TASKS AND IMPLEMENTATION DETAILS

Our experiments study three simulated environments of varying difficulty: 2D point navigation, 2D locomotion using the “HalfCheetah,” and 3D locomotion using the “Ant,” with the latter two environments are modifications of popular reinforcement learning benchmarks (Duan et al., 2016a). While the 2D navigation environment allows for direct control of position, HalfCheetah and Ant can only control their center of mass via feedback control with high dimensional actions (6D for HalfCheetah, 8D for Ant) and observations (17D for HalfCheetah, 111D for Ant).

The evaluation tasks, shown in Figure 6, are similar to prior work (Finn et al., 2017a; Pong et al., 2018): 2D navigation and ant require navigating to goal positions, while the half cheetah must run at different goal velocities. These tasks are not accessible to our algorithm during meta-training. We used the default hyperparameters for MAML across all tasks, varying the meta-batch size according to the number of skills that the discriminator is parameterized by - 50 for pointmass, and 20 for cheetah and ant. We found that the default architecture - 2 layer MLP with 300 units each and ReLU non-linearities worked quite well for meta-training. We also used the default hyperparameters for DIAYN to acquire skills. We swept over learning rates for learning from scratch via vanilla policy gradient, and found that using ADAM with adaptive step size is the most stable and quick at learning.

#### 4.2 FAST ADAPTATION AFTER UNSUPERVISED META LEARNING

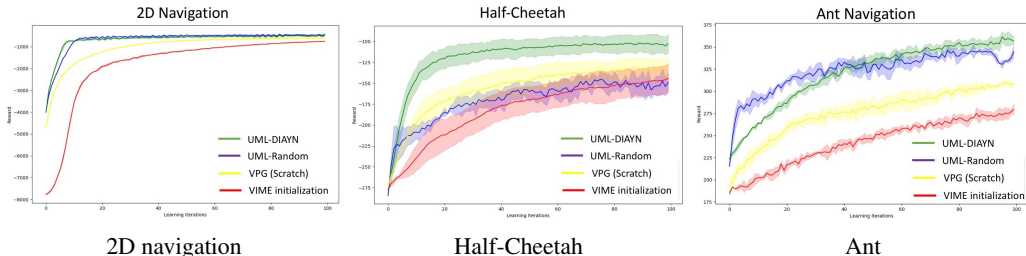


Figure 3: **Unsupervised meta-learning accelerates learning:** After unsupervised meta-learning, our approach (UML-DIAYN and UML-RANDOM) quickly learns a new task significantly faster than learning from scratch, especially on complex tasks. Learning the task distribution with DIAYN helps more for complex tasks. Results are averaged across 20 evaluation tasks, and 3 random seeds for testing. UML-DIAYN and random also significantly outperform learning with DIAYN initialization or an initialization with a policy pretrained with VIME.

The comparison between the two variants of unsupervised meta-learning and learning from scratch is shown in Fig 3, and we compare to hand-crafted task distributions in Fig 4. In Figure 3, we also add a comparison to VIME Houthoof et al. (2016), a standard novelty-based exploration method, where we pretrain a policy with the VIME reward and then finetune it on the meta-test tasks. In all cases, the UML-DIAYN variant of unsupervised meta-learning produces an RL procedure that outperforms reinforcement learning from scratch and VIME-init, suggesting that unsupervised interaction with the environment and meta-learning is effective in producing environment-specific but task-agnostic priors that accelerate learning on new, previously unseen tasks. The comparison with VIME shows that the speed of learning is not just about exploration but is indeed about fast adaptation. In our experiments thus far, UML-DIAYN always performs better than learning from scratch, although the benefit varies across tasks depending on the actual performance of DIAYN.

Interestingly, in many cases (in Fig 4) the performance of unsupervised meta-learning with DIAYN matches that of the hand-designed task distribution. We see that on the 2D navigation task, while handcrafted meta-learning is able to learn very quickly initially, it performs similarly after 100 steps. For the cheetah environment as well, handcrafted meta-learning is able to learn very quickly to start off, but is quickly matched by unsupervised meta-RL with DIAYN. On the ant task, we see that hand-crafted meta-learning does do better than UML-DIAYN, likely because the task distribution is more challenging, and a better unsupervised task proposal algorithm would improve the performance of a meta-learner.

The comparison between the two unsupervised meta-learning variants is also illuminating: while the DIAYN-based variant of our method generally achieves the best performance, even the random discriminator is often able to provide a sufficient diversity of tasks to produce meaningful acceleration over learning from scratch in the case of 2D navigation and ant. This result has two interesting implications. First, it suggests that unsupervised meta-learning is an effective tool for learning an environment prior. Although the performance of unsupervised meta-learning can be improved with better coverage using DIAYN (as seen in Fig 3), even the random discriminator version provides competitive advantages over learning from scratch. Second, the comparison provides a clue for identifying the source of the structure learned through unsupervised meta-learning: though the particular task distribution has an effect on performance, simply interacting with the environment (without structured objectives, using a random discriminator) already allows meta-RL to learn effective adaptation strategies in a given environment. That is, the performance cannot be explained



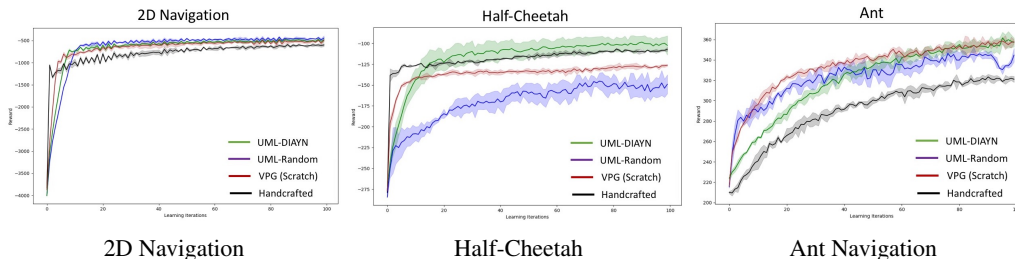


Figure 4: **Comparison with handcrafted tasks:** Unsupervised meta-learning (UML-DIAYN) is competitive with meta-training on handcrafted reward functions (i.e., an oracle). A misspecified, handcrafted meta-training task distribution often performs worse, illustrating the benefits of learning the task distribution.

only by the unsupervised procedure (DIAYN) capturing the right task distribution. We also provide an analysis of the task distributions acquired by the DIAYN procedure in Appendix C.1.

## 5 DISCUSSION AND FUTURE WORK

We presented an unsupervised approach to meta-reinforcement learning, where meta-learning is used to acquire an efficient reinforcement learning procedure without requiring hand-specified task distributions for meta-training. This approach accelerates RL without relying on the manual supervision required for conventional meta-learning algorithms. We provide a theoretical derivation that argues that task proposals based on mutual information maximization can provide for a minimum worst case regret meta-learner under certain assumptions. We then instantiate an approximation to the theoretically motivated method by building on recently developed unsupervised task proposal and meta-learning algorithms. Our experiments indicate that unsupervised meta-RL can accelerate learning on a range of tasks, outperforming learning from scratch and often matching the performance of meta-learning from hand-specified task distributions.

As our work is the first foray into unsupervised meta reinforcement-learning, our approach opens a number of questions about unsupervised meta-learning algorithms. While we focus on purely unsupervised task proposal mechanisms, it is straightforward to incorporate minimally-informative priors into this procedure. For example, we might restrict the learned reward functions to operate on only part of the state. We consider the reinforcement learning setting in our work because environment interaction mediates the unsupervised learning process, ensuring that there is something to learn even without access to task reward.

## REFERENCES

- Joshua Achiam, Harrison Edwards, Dario Amodei, and Pieter Abbeel. Variational option discovery algorithms. *arXiv preprint arXiv:1807.10299*, 2018.
- Maruan Al-Shedivat, Trapti Bansal, Yuri Burda, Ilya Sutskever, Igor Mordatch, and Pieter Abbeel. Continuous adaptation via meta-learning in nonstationary and competitive environments. *arXiv preprint arXiv:1710.03641*, 2017.
- Marcin Andrychowicz, Misha Denil, Sergio Gomez, Matthew W Hoffman, David Pfau, Tom Schaul, and Nando de Freitas. Learning to learn by gradient descent by gradient descent. In *Neural Information Processing Systems (NIPS)*, 2016.
- Marcin Andrychowicz, Filip Wolski, Alex Ray, Jonas Schneider, Rachel Fong, Peter Welinder, Bob McGrew, Josh Tobin, OpenAI Pieter Abbeel, and Wojciech Zaremba. Hindsight experience replay. In *Advances in Neural Information Processing Systems*, pp. 5048–5058, 2017.
- Marc G. Bellemare, Sriram Srinivasan, Georg Ostrovski, Tom Schaul, David Saxton, and Rémi Munos. Unifying count-based exploration and intrinsic motivation. *CoRR*, abs/1606.01868, 2016. URL <http://arxiv.org/abs/1606.01868>.
- Samy Bengio, Yoshua Bengio, Jocelyn Cloutier, and Jan Gecsei. On the optimization of a synaptic learning rule. In *Optimality in Artificial and Biological Neural Networks*, 1992.
- John D Co-Reyes, YuXuan Liu, Abhishek Gupta, Benjamin Eysenbach, Pieter Abbeel, and Sergey Levine. Self-consistent trajectory autoencoder: Hierarchical reinforcement learning with trajectory embeddings. *arXiv preprint arXiv:1806.02813*, 2018.

- Yan Duan, Xi Chen, Rein Houthoofd, John Schulman, and Pieter Abbeel. Benchmarking deep reinforcement learning for continuous control. In *International Conference on Machine Learning*, pp. 1329–1338, 2016a.
- Yan Duan, John Schulman, Xi Chen, Peter L Bartlett, Ilya Sutskever, and Pieter Abbeel.  $RL^2$ : Fast reinforcement learning via slow reinforcement learning. *arXiv preprint arXiv:1611.02779*, 2016b.
- Benjamin Eysenbach, Abhishek Gupta, Julian Ibarz, and Sergey Levine. Diversity is all you need: Learning skills without a reward function. *arXiv preprint arXiv:1802.06070*, 2018.
- Chelsea Finn and Sergey Levine. Meta-learning and universality: Deep representations and gradient descent can approximate any learning algorithm. *CoRR*, abs/1710.11622, 2017. URL <http://arxiv.org/abs/1710.11622>.
- Chelsea Finn and Sergey Levine. Meta-learning and universality: Deep representations and gradient descent can approximate any learning algorithm. *International Conference on Learning Representations*, 2018.
- Chelsea Finn, Pieter Abbeel, and Sergey Levine. Model-agnostic meta-learning for fast adaptation of deep networks. *arXiv preprint arXiv:1703.03400*, 2017a.
- Chelsea Finn, Tianhe Yu, Tianhao Zhang, Pieter Abbeel, and Sergey Levine. One-shot visual imitation learning via meta-learning. *CoRR*, abs/1709.04905, 2017b. URL <http://arxiv.org/abs/1709.04905>.
- Karol Gregor, Danilo Jimenez Rezende, and Daan Wierstra. Variational intrinsic control. *arXiv preprint arXiv:1611.07507*, 2016.
- Abhishek Gupta, Russell Mendonca, YuXuan Liu, Pieter Abbeel, and Sergey Levine. Meta-reinforcement learning of structured exploration strategies. *arXiv preprint arXiv:1802.07245*, 2018.
- David Held, Xinyang Geng, Carlos Florensa, and Pieter Abbeel. Automatic goal generation for reinforcement learning agents. *arXiv preprint arXiv:1705.06366*, 2017.
- Sepp Hochreiter, A Steven Younger, and Peter R Conwell. Learning to learn using gradient descent. In *International Conference on Artificial Neural Networks*, 2001.
- Rein Houthoofd, Xi Chen, Yan Duan, John Schulman, Filip De Turck, and Pieter Abbeel. VIME: variational information maximizing exploration. In *Advances in Neural Information Processing Systems*, 2016.
- Rein Houthoofd, Richard Y Chen, Phillip Isola, Bradly C Stadie, Filip Wolski, Jonathan Ho, and Pieter Abbeel. Evolved policy gradients. *arXiv preprint arXiv:1802.04821*, 2018.
- Philipp Krähenbühl, Carl Doersch, Jeff Donahue, and Trevor Darrell. Data-dependent initializations of convolutional neural networks. *arXiv preprint arXiv:1511.06856*, 2015.
- Ke Li and Jitendra Malik. Learning to optimize. *International Conference on Learning Representations (ICLR)*, 2017.
- Russell Mendonca, Abhishek Gupta, Rosen Kralev, Pieter Abbeel, Sergey Levine, and Chelsea Finn. Guided meta-policy search. *CoRR*, abs/1904.00956, 2019.
- Nikhil Mishra, Mostafa Rohaninejad, Xi Chen, and Pieter Abbeel. A simple neural attentive meta-learner. In *NIPS 2017 Workshop on Meta-Learning*, 2017.
- Tsendsuren Munkhdalai and Hong Yu. Meta networks. *International Conference on Machine Learning (ICML)*, 2017.
- Devang K Naik and RJ Mammone. Meta-neural networks that learn by learning. In *International Joint Conference on Neural Networks (IJCNN)*, 1992.
- Ian Osband, Charles Blundell, Alexander Pritzel, and Benjamin Van Roy. Deep exploration via bootstrapped DQN. *CoRR*, abs/1602.04621, 2016. URL <http://arxiv.org/abs/1602.04621>.
- Deepak Pathak, Pulkit Agrawal, Alexei A. Efros, and Trevor Darrell. Curiosity-driven exploration by self-supervised prediction. In *ICML*, 2017.
- Vitchyr Pong, Shixiang Gu, Murtaza Dalal, and Sergey Levine. Temporal difference models: Model-free deep rl for model-based control. *arXiv preprint arXiv:1802.09081*, 2018.
- Kate Rakelly, Aurick Zhou, Deirdre Quillen, Chelsea Finn, and Sergey Levine. Efficient off-policy meta-reinforcement learning via probabilistic context variables. *arXiv preprint arXiv:1903.08254*, 2019.

- Sachin Ravi and Hugo Larochelle. Optimization as a model for few-shot learning. In *International Conference on Learning Representations (ICLR)*, 2017.
- Jonas Rothfuss, Dennis Lee, Ignasi Clavera, Tamim Asfour, and Pieter Abbeel. Promp: Proximal meta-policy search. In *International Conference on Learning Representations, ICLR*, 2019.
- Adam Santoro, Sergey Bartunov, Matthew Botvinick, Daan Wierstra, and Timothy Lillicrap. Meta-learning with memory-augmented neural networks. In *International Conference on Machine Learning (ICML)*, 2016.
- Tom Schaul, Daniel Horgan, Karol Gregor, and David Silver. Universal value function approximators. In *International Conference on Machine Learning*, pp. 1312–1320, 2015.
- Jürgen Schmidhuber. *Evolutionary principles in self-referential learning, or on learning how to learn: the meta-meta-... hook*. PhD thesis, Technische Universität München, 1987.
- Jürgen Schmidhuber. Driven by compression progress: A simple principle explains essential aspects of subjective beauty, novelty, surprise, interestingness, attention, curiosity, creativity, art, science, music, jokes. In *Computational Creativity: An Interdisciplinary Approach, 12.07. - 17.07.2009*, 2009. URL <http://drops.dagstuhl.de/opus/volltexte/2009/2197/>.
- Jake Snell, Kevin Swersky, and Richard Zemel. Prototypical networks for few-shot learning. In *Advances in Neural Information Processing Systems*, pp. 4080–4090, 2017.
- Bradly C Stadie, Sergey Levine, and Pieter Abbeel. Incentivizing exploration in reinforcement learning with deep predictive models. *arXiv preprint arXiv:1507.00814*, 2015.
- Bradly C. Stadie, Ge Yang, Rein Houthoofd, Xi Chen, Yan Duan, Yuhuai Wu, Pieter Abbeel, and Ilya Sutskever. Some considerations on learning to explore via meta-reinforcement learning. *CoRR*, abs/1803.01118, 2018. URL <http://arxiv.org/abs/1803.01118>.
- Sainbayar Sukhbaatar, Zeming Lin, Ilya Kostrikov, Gabriel Synnaeve, Arthur Szlam, and Rob Fergus. Intrinsic motivation and automatic curricula via asymmetric self-play. *arXiv preprint arXiv:1703.05407*, 2017.
- Flood Sung, Li Zhang, Tao Xiang, Timothy Hospedales, and Yongxin Yang. Learning to learn: Meta-critic networks for sample efficient learning. *arXiv preprint arXiv:1706.09529*, 2017.
- Sebastian Thrun and Lorien Pratt. *Learning to learn*. Springer Science & Business Media, 1998.
- Jane X Wang, Zeb Kurth-Nelson, Dhruva Tirumala, Hubert Soyer, Joel Z Leibo, Remi Munos, Charles Blundell, Dharshan Kumaran, and Matt Botvinick. Learning to reinforcement learn. *arXiv preprint arXiv:1611.05763*, 2016.
- Darrell Whitley and Jean Paul Watson. Complexity theory and the no free lunch theorem, 2005.
- David H Wolpert, William G Macready, et al. No free lunch theorems for search. Technical report, Technical Report SFI-TR-95-02-010, Santa Fe Institute, 1995.

## A APPENDIX

## B PROOFS

**Lemma 1** Let  $\pi$  be a policy for which  $\rho_{\pi}^T(s)$  is uniform. Then  $\pi$  has lowest worst-case regret.

*Proof of Lemma 1.* To begin, we note that all goal distributions  $p(s_g)$  have equal regret for policies where  $\rho_{\pi}^T(s)$  is uniform:

$$\text{REGRET}_p(\pi) = \int \frac{p(s_g)}{\rho_{\pi}^T(s_g)} ds_g = \int \frac{p(s_g)}{|\mathcal{S}|} ds_g = \frac{1}{|\mathcal{S}|}$$

Now, consider a policy  $\pi'$  for which  $\rho_{\pi'}^T(s)$  is not uniform. For simplicity, we will assume that the argmin is unique, though the proof holds for non-unique argmins as well. The worst-case goal distribution will choose the state  $s^-$  where that the policy is least likely to visit:

$$p^-(s_g) \triangleq \mathbb{1}(s_g = \arg \min_s \rho_{\pi'}^T(s))$$

Thus, the worst-case regret for policy  $\pi'$  is strictly greater than the regret for a uniform  $\pi$ :

$$\max_p \text{REGRET}_p(\pi) = \text{REGRET}_{p^-}(\pi) = \int \frac{\mathbb{1}(s_g = \arg \min_s \rho_\pi^T(s))}{\rho_\pi^T(s_g)} ds_g = \frac{1}{\min_s \rho_\pi^T(s)} > \frac{1}{|\mathcal{A}|} \quad (11)$$

Thus, a policy  $\pi'$  for which  $\rho_\pi^T$  is non-uniform cannot be minimax, so the optimal policy has a uniform marginal  $\rho_\pi^T$ .  $\square$

**Lemma 2:** Mutual information  $I(s_T; z)$  is maximized by a task distribution  $p(s_g)$  which is uniform over goal states.

*Proof of Lemma 2.* We define a uniform distribution over goal states  $p(s_g)$  as a latent variable model, where we sample a latent variable  $z$  from a uniform prior  $p(z)$ . In this latent variable model, the marginal entropy  $\mathcal{H}[s_T | z]$  attains the smallest possible value (zero) when each latent variable  $z$  corresponds to exactly one final state,  $s_z$ . In contrast, the marginal entropy  $\mathcal{H}[s_T]$  attains the largest possible value ( $\log |\mathcal{S}|$ ) when the marginal distribution  $p(s_T) = \int p(s_T | z)p(z)dz$  is uniform. Recalling that the mutual information is defined as the difference between these two entropies, we see that maximizing the mutual information results in a uniform distribution over goals.  $\square$

## C ABLATIONS

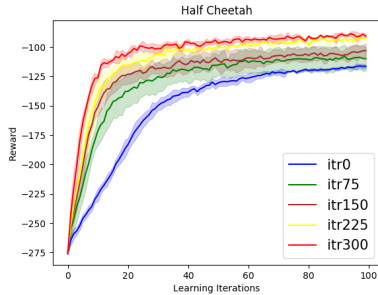


Figure 5: Analysis of effect of additional meta-training on meta-test time learning of new tasks. For larger iterations of meta-trained policies, we have improved test time performance, showing that additional meta-training is beneficial.

To understand the method performance more clearly, we also add an ablation study where we compare the meta-test performance of policies at different iterations along meta-training. This shows the effect that additional meta-training has on the fast learning performance for new tasks. This comparison is shown in Fig 5. As can be seen here, at iteration 0 of meta-training the policy is not a very good initialization for learning new tasks. As we move further along the meta-training process, we see that the meta-learned initialization becomes more and more effective at learning new tasks. This shows a clear correlation between additional meta-training and improved meta test-time performance.

### C.1 ANALYSIS OF LEARNED TASK DISTRIBUTIONS

We can analyze the tasks discovered through unsupervised exploration and compare them to tasks we evaluate on at meta-test time. Figure 6 illustrates these distributions using scatter plots for 2D navigation and the Ant, and a histogram for the HalfCheetah. Note that we visualize dimensions of the state that are relevant for the evaluation tasks – positions and velocities – but these dimensions are *not* specified in any way during unsupervised task acquisition, which operates on the entire state space. Although the tasks proposed via unsupervised exploration provide fairly broad coverage, they are clearly quite distinct from the meta-test tasks, suggesting the approach can tolerate considerable distributional shift. Qualitatively, many of the tasks proposed via unsupervised exploration such as jumping and falling that are not relevant for the evaluation tasks. Our choice of the evaluation tasks was largely based on prior work, and therefore not tailored to this exploration procedure. The results for unsupervised meta-reinforcement learning therefore suggest quite strongly that unsupervised task acquisition can provide an effective meta-training set, at least for MAML, even when evaluating on tasks that do not closely match the discovered task distribution.

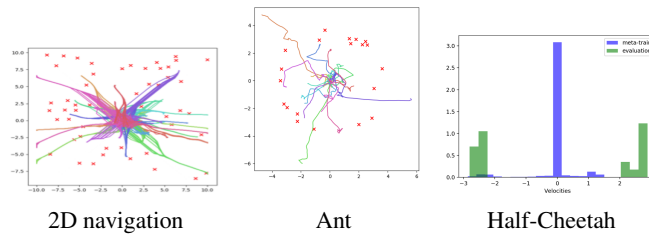


Figure 6: **Learned meta-training task distribution and evaluation tasks:** We plot the center of mass for various skills discovered by point mass and ant using DIAYN, and a blue histogram of goal velocities for cheetah. Evaluation tasks, which are not provided to the algorithm during meta-training, are plotted as red ‘x’ for ant and pointmass, and as a green histogram for cheetah. While the meta-training distribution is broad, it does not fully cover the evaluation tasks. Nonetheless, meta-learning on this *learned* task distribution enables efficient learning on a test task distribution.