
Distributionally Adaptive Meta Reinforcement Learning

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Meta-reinforcement learning algorithms provide a data-driven way to acquire poli-
2 cies that quickly adapt to many tasks with varying rewards or dynamics functions.
3 However, learned meta-policies are often effective only on the exact task distribu-
4 tion on which they were trained and struggle in the presence of distribution shift
5 of test-time rewards or transition dynamics. In this work, we develop a frame-
6 work for meta-RL algorithms that are able to behave appropriately under test-time
7 distribution shifts in the space of tasks. Our framework centers on an adaptive
8 approach to distributional robustness that trains a population of meta-policies to
9 be robust to varying levels of distribution shift. When evaluated on a potentially
10 shifted test-time distribution of tasks, this allows us to choose the meta-policy with
11 the most appropriate level of robustness, and use it to perform fast adaptation. We
12 formally show how our framework allows for improved regret under distribution
13 shift, and empirically show its efficacy on simulated robotics problems under a
14 wide range of distribution shifts.

15 1 Introduction

16 The diversity and dynamism of the real world require reinforcement learning (RL) agents that
17 can quickly adapt and learn new behaviors when placed in novel situations. Meta reinforcement
18 learning provides a framework for conferring this ability to RL agents, by learning a “meta-policy”
19 trained to adapt as quickly as possible to tasks from a provided training distribution [35, 9, 30, 43].
20 Unfortunately, meta-RL agents are prone to overfitting to the distribution of tasks they are trained
21 on, and have been shown to behave erratically when asked to adapt to tasks beyond the training
22 distribution [4, 7]. As an example of this negative transfer, consider using meta-learning to teach
23 a robot to navigate to goals quickly (illustrated in Figure 1). The resulting meta-policy learns to
24 quickly adapt and walk to any target location specified in the training distribution, but explores poorly
25 and fails to adapt to any location not in that distribution. Overfitting is particularly problematic
26 for the meta-learning setting, since the scenarios where we need the ability to learn quickly are
27 usually exactly those where the agent experiences distribution shift. This type of meta-distribution
28 shift afflicts a number of real-world problems including autonomous vehicle driving [8], in-hand
29 manipulation [14, 1], and quadruped locomotion [21, 19, 15], where the test-time task distribution
30 may not be well represented during training.

31 In this work, we study meta-RL algorithms that learn meta-policies resilient to task distribution shift
32 at test time. One approach to enable this resiliency is to leverage the framework of distributional
33 robustness [33], training meta-policies that prepare for distribution shifts by optimizing the *worst-case*
34 empirical risk against a set of task distributions which lie within a bounded distance from the original
35 training task distribution (often referred to as an *uncertainty set*). This allows meta-policies to
36 deal with potential test-time task distribution shift, bounding their worst-case test-time regret for

37 distributional shifts within the chosen uncertainty set. However, choosing an appropriate uncertainty
 38 set can be quite challenging without further information about the test environment, significantly
 39 impacting the test-time performance of algorithms under distribution shift. Large uncertainty sets
 40 allow resiliency to a wider range of distribution shifts, but the resulting meta-policy adapts very slowly
 41 at test time; smaller uncertainty sets enable faster test-time adaptation, but leave the meta-policy
 42 brittle to task distribution shifts. Can we get the best of both worlds?

43 Our key insight is that we can prepare for a variety of
 44 potential test-time distribution shifts by constructing and
 45 training against different uncertainty sets at training time.
 46 By preparing for adaptation against each of these uncer-
 47 tainty sets, an agent is able to adapt to a variety of poten-
 48 tial test-time distribution shifts by adaptively choosing the
 49 most appropriate level of distributional robustness for the
 50 test distribution at hand. We introduce a conceptual frame-
 51 work called distributionally adaptive meta reinforcement
 52 learning formalizing this idea. At train time, the agent
 53 learns robust meta-policies with widening uncertainty sets,
 54 preemptively accounting for different levels of test-time
 55 distribution shift that may be encountered. At test time,
 56 the agent infers the level of distribution shift it is faced
 57 with, and then uses the corresponding meta-policy to adapt
 58 to the new task. In doing so, the agent is able to adaptively
 59 choose the best level of robustness for the test-time task
 60 distribution, preserving the fast adaptation benefits of meta

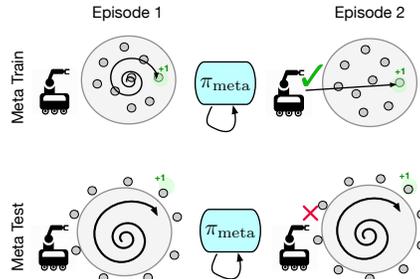


Figure 1: **Failure of Typical Meta-RL.** On meta-training tasks, π_{meta} explores effectively and quickly learns the optimal behavior (top row). When test tasks come from a slightly larger task distribution, exploration fails catastrophically, resulting in poor adaptation behavior (bottom row).

61 RL, while also ensuring good asymptotic performance under distribution shift. We instantiate a
 62 practical algorithm in this framework (DiAMetR), using learned generative models to imagine new
 63 task distributions close to the provided training tasks that can be used to train robust meta-policies.

64 The contribution of this paper is to propose a framework for making meta-reinforcement learning
 65 resilient to a variety of task distribution shifts, and DiAMetR, a practical algorithm instantiating
 66 the framework. DiAMetR trains a population of meta-policies to be robust to different degrees of
 67 distribution shifts and then adaptively chooses a meta-policy to deploy based on the inferred test-time
 68 distribution shift. Our experiments verify the utility of adaptive distributional robustness under
 69 test-time task distribution shift in a number of simulated robotics domains.

70 2 Related Work

71 Meta-reinforcement learning algorithms aim to leverage a distribution of training tasks to “learn a
 72 reinforcement learning algorithm”, that is able to learn as quickly on new tasks drawn from the same
 73 distribution. A variety of algorithms have been proposed for meta-RL, including memory-based
 74 [6, 22], gradient-based [9, 32, 11] and latent-variable based [30, 43, 42] schemes. These algorithms
 75 show the ability to generalize to new tasks drawn from the same distribution, and have been applied
 76 to problems ranging from robotics [24, 42, 15] to computer science education [39]. This line of
 77 work has been extended to operate in scenarios without requiring any pre-specified task distribution
 78 [10, 13] or in offline settings [5, 25, 23] making them more broadly applicable to a wider class of
 79 problems. However, most meta-RL algorithms assume source and target tasks are drawn from the
 80 same distribution, an assumption rarely met in practice. Our work shows how the machinery of
 81 meta-RL can be made compatible with distribution shift at test time, using ideas from distributional
 82 robustness. Some recent work shows that model based meta-reinforcement learning can be made to
 83 be robust to a particular level distribution shift [20, 17] by learning a shared dynamics model against
 84 adversarially chosen task distributions. We show that we can build model-free meta-reinforcement
 85 learning algorithms, which are not just robust to a particular level of distribution shift, but can adapt
 86 to various levels of shift.

87 Distributional robustness methods have been studied in the context of building supervised learning
 88 systems that are robust to the test distribution being different than the training one. The key idea
 89 is to train a model to not just minimize empirical risk, but instead learn a model that has the
 90 lowest worst-case empirical risk among an “uncertainty-set” of distributions that are boundedly close
 91 to the empirical training distribution [33, 18, 2, 12]. If the uncertainty set and optimization are
 92 chosen carefully, these methods have been shown to obtain models that are robust to small amounts

93 of distribution shift at test time [33, 18, 2, 12], finding applications in problems like federated
 94 learning [12] and image classification [18]. This has been extended to the min-max robustness
 95 setting for specific algorithms like model-agnostic meta-learning [3], but are critically dependent on
 96 correct specification of the appropriate uncertainty set and applicable primarily in supervised learning
 97 settings. Alternatively, several RL techniques aim to directly tackle the robustness problem, aiming
 98 to learn policies robust to adversarial perturbations [37, 41, 29, 28]. [40] conditions the policy on
 99 uncertainty sets to make it robust to different perturbation sets. While these methods are able to
 100 learn conservative, robust policies, they are unable to adapt to new tasks as DiAMetR does in the
 101 meta-reinforcement learning setting. In our work, rather than choosing a single uncertainty set, we
 102 learn many meta-policies for widening uncertainty sets, thereby accounting for different levels of
 103 test-time distribution shift.

104 3 Preliminaries

105 **Meta-Reinforcement Learning** aims to learn a fast reinforcement learning algorithm or a “meta-
 106 policy” that can quickly maximize performance on tasks \mathcal{T} from some distribution $p(\mathcal{T})$. Formally,
 107 each task \mathcal{T} is a Markov decision process (MDP) $\mathcal{M} = (\mathcal{S}, \mathcal{A}, \mathcal{P}, \mathcal{R}, \gamma, \mu_0)$; the goal is to exploit
 108 regularities in the structure of rewards and environment dynamics across tasks in $p(\mathcal{T})$ to acquire
 109 effective exploration and adaptation mechanisms that enable learning on new tasks much faster than
 110 learning the task naively from scratch. A meta-policy (or fast learning algorithm) π_{meta} maps a history
 111 of environment experience $h \in (\mathcal{S} \times \mathcal{A} \times \mathcal{R})^*$ in a new task to an action a , and is trained to acquire
 112 optimal behaviors on tasks from $p(\mathcal{T})$ within k episodes:

$$\min_{\pi_{\text{meta}}} \mathbb{E}_{\mathcal{T} \sim p(\mathcal{T})} [\text{Regret}(\pi_{\text{meta}}, \mathcal{T})],$$

$$\text{Regret}(\pi_{\text{meta}}, \mathcal{T}) = J(\pi_{\mathcal{T}}^*) - \mathbb{E}_{a_t^{(i)} \sim \pi_{\text{meta}}(\cdot | h_t^{(i)}), \mathcal{T}} \left[\frac{1}{k} \sum_{i=1}^k \sum_{t=1}^T r_t^{(i)} \right],$$

where $h_t^{(i)} = (s_{1:t}^{(i)}, r_{1:t}^{(i)}, a_{1:t-1}^{(i)}) \cup (s_{1:T}^{(j)}, r_{1:T}^{(j)}, a_{1:T}^{(j)})_{j=1}^{i-1}$. (1)

113 Intuitively, the meta-policy has two components: an exploration mechanism that ensures that appropriate
 114 reward signal is found for all tasks in the training distribution, and an adaptation mechanism
 115 that uses the collected exploratory data to generate optimal actions for the current task. In practice,
 116 the meta-policy may be represented explicitly as an exploration policy conjoined with a policy
 117 update[9, 30], or implicitly as a black-box RNN [6, 43]. We use the terminology “meta-policies”
 118 interchangeably with that of “fast-adaptation” algorithms, since our practical implementation builds
 119 on [27] (which represents the adaptation mechanism using a black-box RNN). Our work focuses
 120 on the setting where there is potential drift between $p_{\text{train}}(\mathcal{T})$, the task distribution we have access to
 121 during training, and $p_{\text{test}}(\mathcal{T})$, the task distribution of interest during evaluation.

122 **Distributional robustness** [33] learns models that do not minimize empirical risk against the training
 123 distribution, but instead prepare for distribution shift by optimizing the *worst-case* empirical risk
 124 against a set of data distributions close to the training distribution (called an *uncertainty set*):

$$\min_{\theta} \max_{\phi} \mathbb{E}_{x \sim q_{\phi}(x)} [l(x; \theta)] \quad \text{s.t.} \quad D(p_{\text{train}}(x) || q_{\phi}(x)) \leq \epsilon \quad (2)$$

125 This optimization finds the model parameters θ that minimizes worst case risk l over distributions
 126 $q_{\phi}(x)$ in an ϵ -ball (measured by an f -divergence) from the training distribution $p_{\text{train}}(x)$.

127 4 Distributionally Adaptive Meta-Reinforcement Learning

128 In this section, we develop a framework for learning meta-policies, that given access to a training
 129 distribution of tasks $p_{\text{train}}(\mathcal{T})$, is still able to adapt to tasks from a test-time distribution $p_{\text{test}}(\mathcal{T})$ that
 130 is similar but not identical to the training distribution. We introduce a framework for distributionally
 131 adaptive meta-RL below and instantiate it as a practical method in Section 5.

132 4.1 Known Level of Test-Time Distribution Shift

133 We begin by studying a simplified problem where we can exactly quantify the degree to which
 134 the test distribution deviates from the training distribution. Suppose we know that p_{test} satisfies

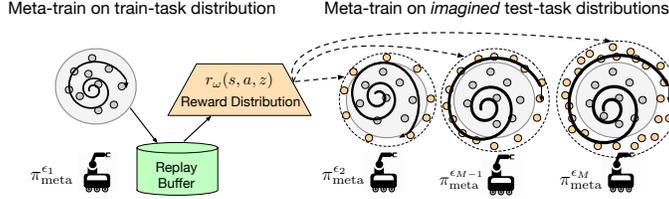


Figure 2: DiAMeTR first learns a meta-policy $\pi_{\text{meta}}^{\epsilon_1}$ and reward distribution $r_\omega(s, a, z)$ on train task distribution. Then, it uses the reward distribution to imagine different shifted test task distributions (orange dots) on which it learns different meta-policies $\{\pi_{\text{meta}}^{\epsilon_i}\}_{i=2}^M$, each corresponding to a different level of robustness .

135 $D(p_{\text{test}}(\mathcal{T})||p_{\text{train}}(\mathcal{T})) < \epsilon$ for some $\epsilon > 0$, where $D(\cdot||\cdot)$ is a probability divergence on the set of task
 136 distributions (e.g. an f -divergence [31] or a Wasserstein distance [36]). A natural learning objective
 137 to learn a meta-policy under this assumption is to minimize the *worst-case* test-time regret across any
 138 test task distribution $q(\mathcal{T})$ that is within some ϵ divergence of the train distribution:

$$\min_{\pi_{\text{meta}}} \mathcal{R}(\pi_{\text{meta}}, p_{\text{train}}(\mathcal{T}), \epsilon),$$

$$\mathcal{R}(\pi_{\text{meta}}, p_{\text{train}}(\mathcal{T}), \epsilon) = \max_{q(\mathcal{T})} \mathbb{E}_{\mathcal{T} \sim q(\mathcal{T})} [\text{Regret}(\pi_{\text{meta}}, \mathcal{T})] \quad \text{s.t. } D(p_{\text{train}}(\mathcal{T})||q(\mathcal{T})) \leq \epsilon \quad (3)$$

139 Solving this optimization problem results in a meta-policy that has been trained to adapt to tasks
 140 from a *wider* task distribution than the original training distribution. It is worthwhile distinguishing
 141 this robust meta-objective, which incentivizes a *robust adaptation mechanism* to a wider set of tasks,
 142 from robust objectives in standard RL, which produce base policies robust to a wider set of dynamics
 143 conditions. The objective in Eq 3 incentivizes an agent to explore and adapt more broadly, not act
 144 more conservatively as standard robust RL methods [29] would encourage. Naturally, the quality of
 145 the robust meta-policy depends on the size of the uncertainty set. If ϵ is large, or the geometry of the
 146 divergence poorly reflect natural task variations, then the robust policy will have to adapt to an overly
 147 large set of tasks, potentially degrading the speed of adaptation.

148 4.2 Handling Arbitrary Levels of Distribution Shift

149 In practice, it is not known how the test distribution p_{test} deviates from the training distribution, and
 150 consequently it is challenging to determine what ϵ to use in the meta-robustness objective. We propose
 151 to overcome this via an adaptive strategy: to train meta-policies for *varying* degrees of distribution
 152 shift, and at test-time, inferring which distribution shift is most appropriate through experience.

153 We train a population of meta-policies $\{\pi_{\text{meta}}^{(i)}\}_{i=1}^M$, each solving the distributionally robust meta-RL
 154 objective (eq 3) for a different level of robustness ϵ_i :

$$\left\{ \pi_{\text{meta}}^{\epsilon_i} := \arg \min_{\pi_{\text{meta}}} \mathcal{R}(\pi_{\text{meta}}, p_{\text{train}}(\mathcal{T}), \epsilon_i) \right\}_{i=1}^M \quad \text{where } \epsilon_M > \epsilon_{M-1} > \dots > \epsilon_1 = 0 \quad (4)$$

155 In choosing a spectrum of ϵ_i , we learn a set of meta-policies that have been trained on increasingly
 156 large set of tasks: at one end ($i = 1$), the meta-policy is trained only on the original training
 157 distribution, and at the other ($i = M$), the meta-policy trained to adapt to any possible task within the
 158 parametric family of tasks. These policies span a tradeoff between being robust to a wider set of task
 159 distributions with larger ϵ (allowing for larger distribution shifts), and being able to adapt quickly to
 160 any given task with smaller ϵ (allowing for better per-task regret minimization).

161 With a set of meta-policies in hand, we must now decide how to leverage test-time experience to
 162 discover the right one to use for the actual test distribution p_{test} . We recognize that the problem
 163 of policy selection can be treated as a stochastic multi-armed bandit problem (precise formulation
 164 in Appendix A), where pulling arm i corresponds to running the meta-policy $\pi_{\text{meta}}^{\epsilon_i}$ for an entire
 165 meta-episode (k task episodes). If a zero-regret bandit algorithm (eg: Thompson’s sampling [38]) is
 166 used, then after a certain number of test-time meta episodes, we can guarantee that the meta-policy
 167 selection mechanism will converge to the meta-policy that best balances the tradeoff between adapting
 168 quickly while still being able to adapt to all the tasks from $p_{\text{test}}(\mathcal{T})$.

169 To summarize our framework for distributionally adaptive meta-RL, we train a population of meta-
 170 policies at varying levels of robustness on a distributionally robust objective that forces the learned

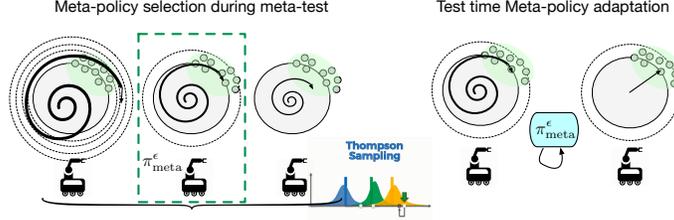


Figure 3: DiAMetR chooses appropriate meta-policy based on inferred distribution shift with Thompson’s sampling and then quickly adapts the selected meta-policy to individual tasks during meta-test.

171 adaptation mechanism to also be robust to tasks not in the training task distribution. At test-time, we
 172 use a bandit algorithm to select the meta-policy whose adaptation mechanism has the best tradeoff
 173 between robustness and speed of adaptation specifically on the test task distribution. Combining
 174 distributional robustness with test-time adaptation allows the adaptation mechanism to work even
 175 if distribution shift is present, while obviating the decreased performance that usually accompanies
 176 overly conservative, distributionally robust solutions.

177 4.3 Analysis

178 To provide some intuition on the properties of this algorithm, we formally analyze adaptive distribu-
 179 tional robustness in a simplified meta RL problem involving tasks \mathcal{T}_g corresponding to reaching some
 180 unknown goal g in a deterministic MDP \mathcal{M} , exactly at the final timestep of an episode. We assume
 181 that all goals are reachable, and use the family of meta-policies that use a stochastic exploratory
 182 policy π until the goal is discovered and return to the discovered goal in all future episodes. The
 183 performance of a meta-policy on a task \mathcal{T}_g under this model can be expressed in terms of the state
 184 distribution of the exploratory policy: $\text{Regret}(\pi_{\text{meta}}, \mathcal{T}_g) = \frac{1}{d_{\pi}^T(g)}$. This particular framework has
 185 been studied in [10, 16], and is a simple, interpretable framework for analysis.

186 We seek to understand performance under distribution shift when the original training task distribution
 187 is relatively concentrated on a subset of possible tasks. We choose the training distribution $p_{\text{train}}(\mathcal{T}_g) =$
 188 $(1 - \beta)\text{Uniform}(\mathcal{S}_0) + \beta\text{Uniform}(\mathcal{S} \setminus \mathcal{S}_0)$, so that p_{train} is concentrated on tasks involving a subset
 189 of the state space $\mathcal{S}_0 \subset \mathcal{S}$, with β a parameter dictating the level of concentration, and consider
 190 test distributions that perturb under the TV metric. Our main result compares the performance of a
 191 meta-policy trained to an ϵ_2 -level of robustness when the true test distribution deviates by ϵ_1 .

192 **Proposition 4.1** *Let $\bar{\epsilon}_i = \min\{\epsilon_i + \beta, 1 - \frac{|\mathcal{S}_0|}{|\mathcal{S}|}\}$. There exists $q(\mathcal{T})$ satisfying $D_{TV}(p_{\text{train}}, q) \leq \epsilon_1$*
 193 *where an ϵ_2 -robust meta policy incurs excess regret over the optimal ϵ_1 -robust meta-policy:*

$$\mathbb{E}_{q(\mathcal{T})}[\text{Regret}(\pi_{\text{meta}}^{\epsilon_1}, \mathcal{T}) - \text{Regret}(\pi_{\text{meta}}^{\epsilon_2}, \mathcal{T})] \geq \left(c(\epsilon_1, \epsilon_2) + \frac{1}{c(\epsilon_1, \epsilon_2)} - 2 \right) \sqrt{\epsilon_1(1 - \bar{\epsilon}_1)|\mathcal{S}_0|(|\mathcal{S}| - |\mathcal{S}_0|)} \quad (5)$$

194 *The scale of regret depends on $c(\epsilon_1, \epsilon_2) = \sqrt{\frac{\epsilon_2^{-1} - 1}{\epsilon_1^{-1} - 1}}$, a measure of the mismatch between ϵ_1 and ϵ_2 .*

195 We first compare robust and non-robust solutions by analyzing the bound when $\epsilon_2 = 0$. In the regime
 196 of $\beta \ll 1$, excess regret scales as $\mathcal{O}(\epsilon_1 \sqrt{\frac{1}{\beta}})$, meaning that the robust solution is most necessary
 197 when the training distribution is highly concentrated in a subset of the task space. At one extreme, if
 198 the training distribution contains no examples of tasks outside \mathcal{S}_0 ($\beta = 0$), the non-robust solution
 199 incurs *infinite excess regret*; at the other extreme, if the training distribution is uniform on the set of
 200 all possible tasks ($\beta = 1 - \frac{|\mathcal{S}_0|}{|\mathcal{S}|}$), *robustness provides no benefit*.

201 We next quantify the effect of mis-specifying the level of robustness in the meta-robustness objective,
 202 and what benefits *adaptive* distributional robustness can confer. For small β and fixed ϵ_1 , the excess
 203 regret of an ϵ_2 -robust policy scales as $\mathcal{O}(\sqrt{\max\{\frac{\epsilon_2}{\epsilon_1}, \frac{\epsilon_1}{\epsilon_2}\}})$, meaning that excess regret gets incurred if
 204 the meta-policy is trained either to be too robust ($\epsilon_2 \gg \epsilon_1$) or not robust enough $\epsilon_1 \gg \epsilon_2$. Compared
 205 to a fixed robustness level, our strategy of training meta-policies for a sequence of robustness levels

Algorithm 1 DiAMetR: Meta-training phase

- 1: Given: $p_{\text{train}}(\mathcal{T})$, Return: Π
- 2: $\pi_{\text{meta},\theta}^{\epsilon_1}, \mathcal{D}_{\text{Replay-Buffer}} \leftarrow$ Solve equation 1 with off-policy RL²
- 3: Reward distribution r_ω , prior $p_{\text{train}}(z) \leftarrow$ Solve eq 7 using $\mathcal{D}_{\text{Replay-Buffer}}$
- 4: **for** ϵ in $\{\epsilon_2, \dots, \epsilon_M\}$ **do**
- 5: Initialize $q_\phi(z)$, $\pi_{\text{meta},\theta}^\epsilon$ and $\lambda \geq 0$.
- 6: **for** iteration $n = 1, 2, \dots$ **do**
- 7: **Meta-policy:** Update $\pi_{\text{meta},\theta}^\epsilon$ using off-policy RL² [27]

$$\theta := \theta + \alpha \nabla_\theta \mathbb{E}_{z \sim q_\phi(z)} (\mathbb{E}_{\pi_{\text{meta},\theta}^\epsilon, \mathcal{P}} (\frac{1}{k} \sum_{i=1}^k \sum_{t=1}^T r_\omega(s_t^{(i)}, a_t^{(i)}, z)))$$

- 8: **Adversarial task distribution:** Update q_ϕ using Reinforce [34]

$$\phi := \phi - \alpha \nabla_\phi (\mathbb{E}_{z \sim q_\phi(z)} [\mathbb{E}_{\pi_{\text{meta},\theta}^\epsilon, \mathcal{P}} [\frac{1}{k} \sum_{i=1}^k \sum_{t=1}^T r_\omega(s_t^{(i)}, a_t^{(i)}, z)]] + \lambda D_{\text{KL}}(p_{\text{train}}(z) \| q_\phi(z)))$$

- 9: **Lagrange constraint multiplier:** Update λ to enforce $D_{\text{KL}}(p_{\text{train}}(z) \| q_\phi(z)) < \epsilon$,

$$\lambda := \lambda_{\geq 0} + \alpha (D_{\text{KL}}(p_{\text{train}}(z) \| q_\phi(z)) - \epsilon)$$

- 10: **end for**
 - 11: **end for**
-

206 $\{\epsilon_i\}_{i=1}^M$ ensures that this misspecification constant is at most the relative spacing between robustness
207 levels: $\max_i \frac{\epsilon_i}{\epsilon_{i-1}}$. This enables the distributionally adaptive approach to *control* the amount of excess
208 regret by making the sequence more fine-grained, while a fixed choice of robustness incurs larger
209 regret (as we verify empirically in our experiments as well).

210 5 DiAMetR: A Practical Algorithm for Meta-Distribution Shift

211 In order to instantiate our distributionally adaptive framework into a practical algorithm, we must
212 address how task distributions should be parameterized and optimized over, and also how the robust
213 meta-RL problem can be solved with stochastic gradient methods. For simplicity, in the remainder
214 of the paper, we focus on the setting where tasks share transition dynamics, but have different
215 reward functions. We first introduce the individual components of task parameterization and robust
216 optimization, and describe the overall algorithm in Algorithm 1 and 2.

217 **Parameterizing Task Distributions:** Since we assume that variations in tasks correspond to changes
218 in the reward function, the problem of representing a task distribution reduces to learning distributions
219 over reward functions. We propose to learn a probabilistic model of the task reward functions seen
220 in the training task distribution, and use the learned latent representation as a space on which to
221 parameterize uncertainty sets over new task distributions. Specifically, we jointly train a reward
222 encoder $q_\psi(z|h)$ that encodes reward samples from an environment history into the latent space, and
223 a decoder $r_\omega(s, a, z)$ mapping a latent vector z to a reward function using a dataset of trajectories
224 collected from the training tasks. This generative model over reward functions can be trained as a
225 standard latent variable model by maximizing a standard evidence lower bound (ELBO), trading off
226 reward prediction and matching a prior $p_{\text{train}}(z)$ (chosen to be the unit gaussian).

$$\min_{\omega, \psi} \mathbb{E}_{h \sim \mathcal{D}} \left[\mathbb{E}_{z \sim q_\psi(z|h)} \left[\sum_{t=1}^T (r_\omega(s_t, a_t, z) - r_t)^2 \right] + D_{\text{KL}}(q_\psi(z|h) \| \mathcal{N}(0, I)) \right] \quad (6)$$

227 Having learned a latent space, we can parameterize new task distributions $q(\mathcal{T})$ as distributions $q_\phi(z)$
228 (the original training distribution corresponds to $p_{\text{train}}(z) = \mathcal{N}(0, I)$, and measure the divergence
229 between task distributions as well using the KL divergence in this latent space $D(p_{\text{train}}(z) \| q_\phi(z))$).

230 **Learning Robust Meta-Policies:** Given this task parameterization, the next question becomes how to
231 actually perform the robust optimization laid out in Eq:3. The distributional meta-robustness objective
232 can be modelled as an adversarial game between a meta-policy $\pi_{\text{meta}}^\epsilon$ and a task proposal distribution
233 $q(\mathcal{T})$. As described above, this task proposal distribution is parameterized as a distribution over latent

Environment	Task reward	r_{train}	$\{r_{\text{test}}^i\}_{i=1}^K$	Θ
*-navigation	$1[\ \text{agent} - \text{target} \ _2 \leq \delta]$	0.50	$\{0.55, 0.60, 0.65, 0.70\}$	2π
Fetch reach	$1[\ \text{gripper} - \text{target} \ _2 \leq \delta]$	0.10	$\{0.12, 0.14, 0.16, 0.18, 0.20\}$	2π
Blocker push	$1[\ \text{block} - \text{target} \ _2 \leq \delta]$	0.50	$\{0.60, 0.70, 0.80, 0.90, 1.0\}$	$\pi/2$

Table 1: Parameters for train task distribution $p_{\text{train}}(s_t) = \{(\Delta \cos \theta, \Delta \sin \theta) \mid \Delta \sim \mathcal{U}(0, r_{\text{train}}), \theta \sim \mathcal{U}(0, \Theta)\}$ and test task distributions $\{p_{\text{test}}^i(s_t) = \{(\Delta \cos \theta, \Delta \sin \theta) \mid \Delta \sim \mathcal{U}(r_{\text{test}}^{i-1}, r_{\text{test}}^i), \theta \sim \mathcal{U}(0, \Theta)\}\}_{i=1}^K$ (where $r_{\text{test}}^0 = r_{\text{train}}$) for different environments

space $q_\phi(z)$, while $\pi_{\text{meta}}^\epsilon$ is parameterized a typical recurrent neural network policy as in [27]. We parameterize $\{\pi_{\text{meta}}^{\epsilon_i}\}_{i=1}^M$ as a discrete set of meta-policies, with one for each chosen value of ϵ .

This leads to a simple alternating optimization scheme (see Algorithm 1), where the meta-policy is trained using a standard meta-RL algorithm (we use off-policy RL² [27] as a base learner), and the task proposal distribution with a constrained optimization method (we use dual gradient descent [26]). Each iteration n , three updates are performed: 1) the meta-policy π_{meta} updated to improve performance on the current task distribution, 2) the task distribution $q(z)$ updated to increase weight on tasks where the current meta-policy adapts poorly and decreases weight on tasks that the current meta-policy can learn, while staying close to the original training distribution, and 3) a penalty coefficient λ is updated to ensure that $q(z)$ satisfies the divergence constraint.

Test-time meta-policy selection: Since test-time meta-policy selection can be framed as a multi-armed bandit problem, we use a generic Thompson’s sampling [38] algorithm (see Algorithm 2). Each meta-episode n , we sample a meta-policy $\pi_{\text{meta}}^\epsilon$ with probability proportional to its estimated average episodic reward, run the sampled meta-policy $\pi_{\text{meta}}^\epsilon$ for an meta-episode (k environment episodes) and then update the estimate of the average episodic reward. Since Thompson’s sampling is a zero-regret bandit algorithm, it will converge to the meta-policy that achieves the highest average episodic reward and lowest regret on the test task distribution.

Algorithm 2 DiAMetR: Meta-test phase

- 1: Given: $p_{\text{test}}(\mathcal{T}), \Pi = \{\pi_{\text{meta}, \theta}^{\epsilon_i}\}_{i=1}^M$
 - 2: Initialize TS = Thompson-Sampler()
 - 3: **for** meta-episode $n = 1, 2, \dots$ **do**
 - 4: Choose meta-policy $i = \text{TS.sample}()$
 - 5: Run $\pi_{\text{meta}, \theta}^{\epsilon_i}$ for meta-episode
 - 6: TS.update(
 arm= i ,
 reward=meta-episode return)
 - 7: **end for**
-

6 Experimental Evaluation

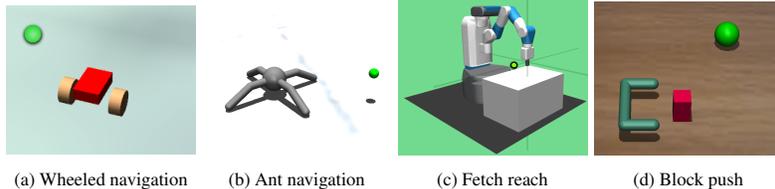


Figure 4: The agent needs to either navigate, move its gripper or push the block to an unobserved target location, indicated by green sphere, by exploring its environment and experiencing reward.

We aim to comprehensively evaluate DiAMetR and answer the following questions: (1) Do meta-policies learned via DiAMetR allow for quick adaptation under different distribution shifts in the test-time task distribution? (2) Does learning for multiple levels of robustness actually help the algorithm adapt more effectively than a particular chosen uncertainty level? (3) Does proposing uncertainty sets via generative modeling provide useful distributions of tasks for robustness?

Setup. We train DiAMetR on four continuous control environments: Wheeled navigation [11] (Wheeled driving a differential drive robot), Ant navigation (Ant controlling a four legged robotic quadruped), Fetch reach and Block push [11] (Figures 4a to 4d) (see Appendix H for more details). Each environment has a train task distribution $\mathcal{T}_i \sim p_{\text{train}}(\mathcal{T})$ such that each task \mathcal{T}_i parameterizes a reward function $r_i(s, a) := r(s, a, \mathcal{T}_i)$. \mathcal{T}_i itself remains unobserved, the agent

268 simply has access to reward values and executing actions in the environment. The meta-policies
 269 are evaluated on train task distribution $p_{\text{train}}(\mathcal{T})$ and on different distributionally shifted test task
 270 distribution $\{p_{\text{test}}^i(\mathcal{T})\}_{i=1}^K$. We use 4 random seeds for all our experiments and include the standard
 271 error bars in our plots. In all of these problems, the distribution of train and test tasks is determined
 272 by the distribution of the underlying target locations s_t , which determines the reward function (exact
 273 distributions in Table 1). Since these environments have sparse rewards, DiAMetR uses a structured
 274 VAE to model reward distributions (see Appendix C for more details).

275 6.1 Adaptation to Varying Levels of Distribution Shift

276 During meta test, given a test task distribution $p_{\text{test}}(\mathcal{T})$, DiAMetR uses Thompson sampling to select
 277 the appropriate meta-policy $\pi_{\text{meta},\theta}^\epsilon$ within $N = 250$ meta episodes. $\pi_{\text{meta},\theta}^\epsilon$ can then solve any task
 278 $\mathcal{T} \sim p_{\text{test}}(\mathcal{T})$ within 1 meta episode ($k = 2$ environment episode). Since DiAMetR adaptively
 279 chooses a meta-policy during test time, we compare it to RL² with test time finetuning. Figure 5
 280 shows that RL²'s performance more or less remains the same after test time finetuning showing
 281 that 10 iteration (with 25 meta-episodes per iteration) isn't enough for RL² to learn a meta-policy
 282 for a new task distribution. For comparison, RL² takes 1500 iterations (with 25 meta-episodes per
 283 iteration) during training to learn a meta-policy for train task distribution.

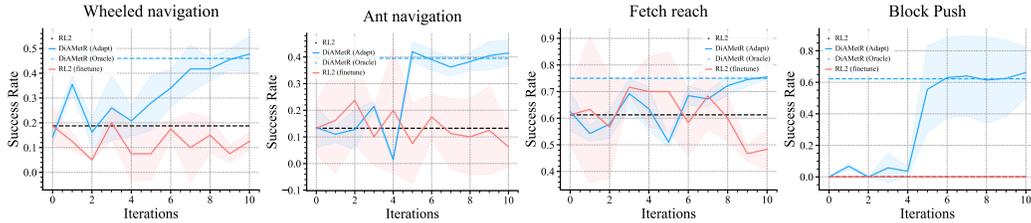


Figure 5: We compare test time adaptation of DiAMetR with test time finetuning of RL² on different environments. We run the adaptation procedure for 10 iterations collecting 25 meta-episodes per iteration. The test target distance distribution for {Wheeled, Ant}-navigation is $\mathcal{U}(0.65, 0.70)$, for Fetch reach is $\mathcal{U}(0.65, 0.70)$ and for Block push is $\mathcal{U}(0.9, 1.0)$. We provide test time adaptation comparisons on other test target distance distributions in Appendix G.

284 To test DiAMetR's ability to adapt to varying levels of distribution shift, we evaluate it on the
 285 above mentioned test task distributions. We compare DiAMetR with meta RL algorithms such as
 286 (off-policy) RL² [27], VariBAD [43] and HyperX [44]. Figure 6 shows that DiAMetR outperforms
 287 RL², VariBAD and HyperX on test task distributions. Furthermore, the performance gap between
 288 DiAMetR and other baselines increase as distribution shift between test task distribution and train task
 289 distribution increases. Naturally, the performance of DiAMetR also deteriorates as the distribution
 290 shift is increased, but as shown in Fig 6, it does so much more slowly than other algorithms. We
 291 also evaluate DiAMetR on train task distribution to see if it incurs any performance loss. Figure 6
 292 shows that DiAMetR either matches or outperforms RL², VariBAD, and HyperX on the train task
 293 distribution. We refer readers to Appendix D for results on point-navigation environment and
 294 Appendix E for ablation studies and further experimental evaluations.

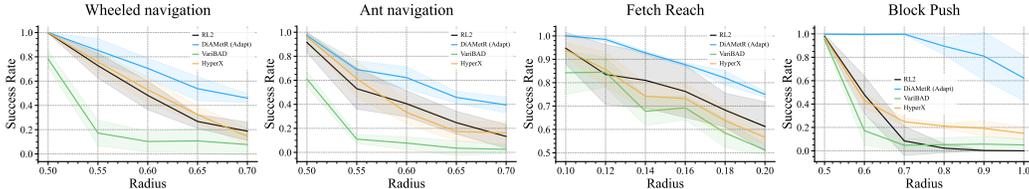


Figure 6: We evaluate DiAMetR and meta RL algorithms (RL², VariBAD and HyperX) on training task distribution and different test task distributions. DiAMetR outperforms RL², VariBAD and HyperX on train distributions and different test distributions. The first point r_{train} on the horizontal axis indicates the training target distance Δ distribution $\mathcal{U}(0, r_{\text{train}})$ and the subsequent points r_{test}^i indicate the shifted test target Δ distribution $\mathcal{U}(r_{\text{test}}^{i-1}, r_{\text{test}}^i)$.

295 **6.2 Analysis of Tasks Proposed by Latent Conditional Uncertainty Sets**

296 We visualize the imagined test reward distribu-
 297 tion for various distribution shifts. Specifically,
 298 we create a heatmap of imagined test reward
 299 functions. Figure 7 visualizes the imagined test
 300 reward distribution in Ant-navigation environ-
 301 ment in increasing order of distribution shifts
 302 with respect to train reward distribution (with
 303 distribution shift parameter ϵ increasing from
 304 left to right). The train distribution of rewards
 305 has uniformly distributed target locations within
 306 the red circle. As clearly seen in Figure 7, as
 307 we increase the distribution shifts, the learned
 308 reward distribution model imagines more target
 309 locations outside the red circle.

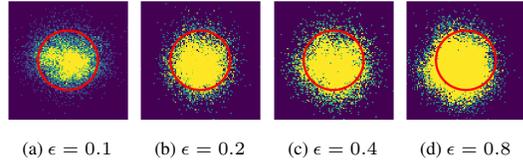


Figure 7: Imagined test reward distributions in Ant-navigation environment in increasing order of distribution shifts. Train reward distribution is uniform within the red circle.

310 **6.3 Analysis of Importance of Multiple Uncertainty Sets**

311 DiAMetR meta-learns a family of adaptation policies, each conditioned on different uncertainty set.
 312 As discussed in section 4, selecting a policy conditioned on a large uncertainty set would lead to
 313 overly conservative behavior. Furthermore, selecting a policy conditioned on a small uncertainty set
 314 would result in failure if the test time distribution shift is high. Therefore, we need to adaptively
 315 select an uncertainty set during test time. To validate this phenomenon empirically, we performed
 316 an ablation study in Figure 8. As clearly visible, adaptively choosing an uncertainty set during test
 317 time allows for better test time distribution adaptation when compared to selecting an uncertainty
 318 set beforehand or selecting a large uncertainty set. These results suggest that a combination of
 319 training robust meta-learners and constructing various uncertainty sets allows for effective test-time
 320 adaptation under distribution shift. DiAMetR is able to avoid both overly conservative behavior and
 321 under-exploration at test-time.

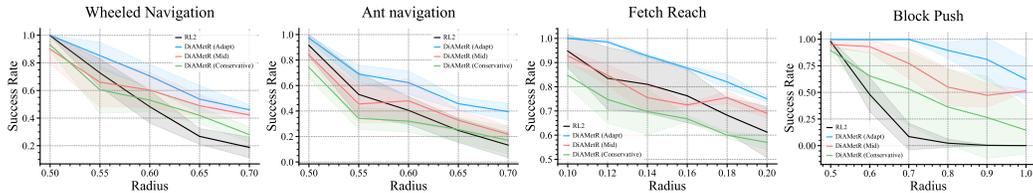


Figure 8: Adaptively choosing an uncertainty set for DiAMetR policy (Adapt) during test time allows it to better adapt to test time distribution shift than choosing an uncertainty set beforehand (Mid). Choosing a large uncertainty set for DiAMetR policy (Conservative) leads to a conservative behavior and hurts its performance when test time distribution shift is low. The first point r_{train} on the horizontal axis indicates the training target distance Δ distribution $\mathcal{U}(0, r_{\text{train}})$ and the subsequent points r_{test}^i indicate the shifted test target distance Δ distribution $\mathcal{U}(r_{\text{test}}^{i-1}, r_{\text{test}}^i)$.

322 **7 Discussion**

323 In this work, we discussed the challenge of distribution shift in meta-reinforcement learning and
 324 showed how we can build meta-reinforcement learning algorithms that are robust to varying levels
 325 of distribution shift. We show how we can build distributionally “adaptive” reinforcement learning
 326 algorithms that can adapt to varying levels of distribution shift, retaining a tradeoff between fast
 327 learning and maintaining asymptotic performance. We then show we can instantiate this algorithm
 328 practically by parameterizing uncertainty sets with a learned generative model. We empirically
 329 showed that this allows for learning meta-learners robust to changes in task distribution.

330 There are several avenues for future work we are keen on exploring, for instance extending adaptive
 331 distributional robustness to more complex meta RL tasks, including those with differing transition
 332 dynamics. Another interesting direction would be to develop a more formal theory providing adaptive
 333 robustness guarantees in meta-RL problems under these inherent distribution shifts.

References

- 334
- 335 [1] T. Chen, J. Xu, and P. Agrawal. A system for general in-hand object re-orientation. In A. Faust,
336 D. Hsu, and G. Neumann, editors, *Proceedings of the 5th Conference on Robot Learning*,
337 volume 164 of *Proceedings of Machine Learning Research*, pages 297–307. PMLR, 08–11 Nov
338 2022. URL <https://proceedings.mlr.press/v164/chen22a.html>.
- 339 [2] J. Cohen, E. Rosenfeld, and Z. Kolter. Certified adversarial robustness via randomized smooth-
340 ing. In *International Conference on Machine Learning*, pages 1310–1320. PMLR, 2019.
- 341 [3] L. Collins, A. Mokhtari, and S. Shakkottai. Distribution-agnostic model-agnostic meta-learning.
342 *CoRR*, abs/2002.04766, 2020. URL <https://arxiv.org/abs/2002.04766>.
- 343 [4] T. Deleu and Y. Bengio. The effects of negative adaptation in model-agnostic meta-learning.
344 *arXiv preprint arXiv:1812.02159*, 2018.
- 345 [5] R. Dorfman, I. Shenfeld, and A. Tamar. Offline meta learning of exploration. *arXiv preprint*
346 *arXiv:2008.02598*, 2020.
- 347 [6] Y. Duan, J. Schulman, X. Chen, P. L. Bartlett, I. Sutskever, and P. Abbeel. RL2: Fast reinforce-
348 ment learning via slow reinforcement learning. *arXiv preprint arXiv:1611.02779*, 2016.
- 349 [7] A. Fallah, A. Mokhtari, and A. Ozdaglar. Generalization of model-agnostic meta-learning
350 algorithms: Recurring and unseen tasks. *Advances in Neural Information Processing Systems*,
351 34, 2021.
- 352 [8] A. Filos, P. Tigkas, R. McAllister, N. Rhinehart, S. Levine, and Y. Gal. Can autonomous
353 vehicles identify, recover from, and adapt to distribution shifts? In *International Conference on*
354 *Machine Learning*, pages 3145–3153. PMLR, 2020.
- 355 [9] C. Finn, P. Abbeel, and S. Levine. Model-agnostic meta-learning for fast adaptation of deep
356 networks. In *International conference on machine learning*, pages 1126–1135. PMLR, 2017.
- 357 [10] A. Gupta, B. Eysenbach, C. Finn, and S. Levine. Unsupervised meta-learning for reinforcement
358 learning. *arXiv preprint arXiv:1806.04640*, 2018.
- 359 [11] A. Gupta, R. Mendonca, Y. Liu, P. Abbeel, and S. Levine. Meta-reinforcement learning of
360 structured exploration strategies. *Advances in neural information processing systems*, 31, 2018.
- 361 [12] J. Hong, H. Wang, Z. Wang, and J. Zhou. Federated robustness propagation: Sharing adversarial
362 robustness in federated learning. *arXiv preprint arXiv:2106.10196*, 2021.
- 363 [13] A. Jabri, K. Hsu, A. Gupta, B. Eysenbach, S. Levine, and C. Finn. Unsupervised curricula for
364 visual meta-reinforcement learning. *Advances in Neural Information Processing Systems*, 32,
365 2019.
- 366 [14] L. Ke, J. Wang, T. Bhattacharjee, B. Boots, and S. Srinivasa. Grasping with chopsticks:
367 Combating covariate shift in model-free imitation learning for fine manipulation. In *2021 IEEE*
368 *International Conference on Robotics and Automation (ICRA)*, pages 6185–6191. IEEE, 2021.
- 369 [15] A. Kumar, Z. Fu, D. Pathak, and J. Malik. Rma: Rapid motor adaptation for legged robots.
370 *arXiv preprint arXiv:2107.04034*, 2021.
- 371 [16] L. Lee, B. Eysenbach, E. Parisotto, E. P. Xing, S. Levine, and R. Salakhutdinov. Efficient
372 exploration via state marginal matching. *CoRR*, abs/1906.05274, 2019. URL <http://arxiv.org/abs/1906.05274>.
- 373
- 374 [17] Z. Lin, G. Thomas, G. Yang, and T. Ma. Model-based adversarial meta-reinforcement
375 learning. In H. Larochelle, M. Ranzato, R. Hadsell, M. Balcan, and H. Lin, ed-
376 itors, *Advances in Neural Information Processing Systems 33: Annual Conference*
377 *on Neural Information Processing Systems 2020, NeurIPS 2020, December 6-12,*
378 *2020, virtual*, 2020. URL <https://proceedings.neurips.cc/paper/2020/hash/73634c1dcbe056c1f7dcf5969da406c8-Abstract.html>.
- 379

- 380 [18] A. Madry, A. Makelov, L. Schmidt, D. Tsipras, and A. Vladu. Towards deep learning models
381 resistant to adversarial attacks. *arXiv preprint arXiv:1706.06083*, 2017.
- 382 [19] G. B. Margolis, G. Yang, K. Paigwar, T. Chen, and P. Agrawal. Rapid locomotion via reinforce-
383 ment learning. *arXiv preprint arXiv:2205.02824*, 2022.
- 384 [20] R. Mendonca, X. Geng, C. Finn, and S. Levine. Meta-reinforcement learning robust to distribu-
385 tional shift via model identification and experience relabeling. *CoRR*, abs/2006.07178, 2020.
386 URL <https://arxiv.org/abs/2006.07178>.
- 387 [21] T. Miki, J. Lee, J. Hwangbo, L. Wellhausen, V. Koltun, and M. Hutter. Learning robust
388 perceptive locomotion for quadrupedal robots in the wild. *Science Robotics*, 7(62):eabk2822,
389 2022.
- 390 [22] N. Mishra, M. Rohaninejad, X. Chen, and P. Abbeel. A simple neural attentive meta-learner.
391 *arXiv preprint arXiv:1707.03141*, 2017.
- 392 [23] E. Mitchell, R. Rafailov, X. B. Peng, S. Levine, and C. Finn. Offline meta-reinforcement
393 learning with advantage weighting. In *International Conference on Machine Learning*, pages
394 7780–7791. PMLR, 2021.
- 395 [24] A. Nagabandi, I. Clavera, S. Liu, R. S. Fearing, P. Abbeel, S. Levine, and C. Finn. Learning to
396 adapt in dynamic, real-world environments through meta-reinforcement learning. *arXiv preprint*
397 *arXiv:1803.11347*, 2018.
- 398 [25] A. Nair, A. Gupta, M. Dalal, and S. Levine. Awac: Accelerating online reinforcement learning
399 with offline datasets. *arXiv preprint arXiv:2006.09359*, 2020.
- 400 [26] Y. Nesterov. Primal-dual subgradient methods for convex problems. *Mathematical programming*,
401 120(1):221–259, 2009.
- 402 [27] T. Ni, B. Eysenbach, S. Levine, and R. Salakhutdinov. Recurrent model-free RL is a strong base-
403 line for many POMDPs, 2022. URL <https://openreview.net/forum?id=E0zOKxQsZhN>.
- 404 [28] T. P. Oikarinen, W. Zhang, A. Megretski, L. Daniel, and T. Weng. Robust deep reinforcement
405 learning through adversarial loss. In M. Ranzato, A. Beygelzimer, Y. N. Dauphin, P. Liang,
406 and J. W. Vaughan, editors, *Advances in Neural Information Processing Systems 34: Annual*
407 *Conference on Neural Information Processing Systems 2021, NeurIPS 2021, December 6-*
408 *14, 2021, virtual*, pages 26156–26167, 2021. URL [https://proceedings.neurips.cc/](https://proceedings.neurips.cc/paper/2021/hash/dbb422937d7ff56e049d61da730b3e11-Abstract.html)
409 [paper/2021/hash/dbb422937d7ff56e049d61da730b3e11-Abstract.html](https://proceedings.neurips.cc/paper/2021/hash/dbb422937d7ff56e049d61da730b3e11-Abstract.html).
- 410 [29] L. Pinto, J. Davidson, R. Sukthankar, and A. Gupta. Robust adversarial reinforcement learning.
411 In D. Precup and Y. W. Teh, editors, *Proceedings of the 34th International Conference on*
412 *Machine Learning, ICML 2017, Sydney, NSW, Australia, 6-11 August 2017*, volume 70 of
413 *Proceedings of Machine Learning Research*, pages 2817–2826. PMLR, 2017. URL <http://proceedings.mlr.press/v70/pinto17a.html>.
414
- 415 [30] K. Rakelly, A. Zhou, C. Finn, S. Levine, and D. Quillen. Efficient off-policy meta-reinforcement
416 learning via probabilistic context variables. In *International conference on machine learning*,
417 pages 5331–5340. PMLR, 2019.
- 418 [31] A. Rényi. On measures of entropy and information. In *Proceedings of the Fourth Berkeley*
419 *Symposium on Mathematical Statistics and Probability, Volume 1: Contributions to the Theory*
420 *of Statistics*, volume 4, pages 547–562. University of California Press, 1961.
- 421 [32] J. Rothfuss, D. Lee, I. Clavera, T. Asfour, and P. Abbeel. Promp: Proximal meta-policy search.
422 *arXiv preprint arXiv:1810.06784*, 2018.
- 423 [33] A. Sinha, H. Namkoong, R. Volpi, and J. Duchi. Certifying some distributional robustness with
424 principled adversarial training. *arXiv preprint arXiv:1710.10571*, 2017.
- 425 [34] R. S. Sutton, D. McAllester, S. Singh, and Y. Mansour. Policy gradient methods for reinforce-
426 ment learning with function approximation. *Advances in neural information processing systems*,
427 12, 1999.

- 428 [35] S. Thrun and L. Y. Pratt, editors. *Learning to Learn*. Springer, 1998. ISBN 978-
 429 1-4613-7527-2. doi: 10.1007/978-1-4615-5529-2. URL [https://doi.org/10.1007/](https://doi.org/10.1007/978-1-4615-5529-2)
 430 [978-1-4615-5529-2](https://doi.org/10.1007/978-1-4615-5529-2).
- 431 [36] L. N. Vaserstein. Markov processes over denumerable products of spaces, describing large
 432 systems of automata. *Problemy Peredachi Informatsii*, 5(3):64–72, 1969.
- 433 [37] E. Vinitzky, Y. Du, K. Parvate, K. Jang, P. Abbeel, and A. M. Bayen. Robust reinforcement
 434 learning using adversarial populations. *CoRR*, abs/2008.01825, 2020. URL [https://arxiv.](https://arxiv.org/abs/2008.01825)
 435 [org/abs/2008.01825](https://arxiv.org/abs/2008.01825).
- 436 [38] D. Wolpert and W. Macready. No free lunch theorems for optimization. *IEEE Transactions on*
 437 *Evolutionary Computation*, 1(1):67–82, 1997. doi: 10.1109/4235.585893.
- 438 [39] M. Wu, N. Goodman, C. Piech, and C. Finn. Prototransformer: A meta-learning approach to
 439 providing student feedback. *arXiv preprint arXiv:2107.14035*, 2021.
- 440 [40] A. Xie, S. Sodhani, C. Finn, J. Pineau, and A. Zhang. Robust policy learning over multiple
 441 uncertainty sets. *arXiv preprint arXiv:2202.07013*, 2022.
- 442 [41] H. Zhang, H. Chen, D. S. Boning, and C. Hsieh. Robust reinforcement learning on state
 443 observations with learned optimal adversary. In *9th International Conference on Learning*
 444 *Representations, ICLR 2021, Virtual Event, Austria, May 3-7, 2021*. OpenReview.net, 2021.
 445 URL <https://openreview.net/forum?id=sCZbhBvqQaU>.
- 446 [42] T. Z. Zhao, A. Nagabandi, K. Rakelly, C. Finn, and S. Levine. Meld: Meta-reinforcement
 447 learning from images via latent state models. *arXiv preprint arXiv:2010.13957*, 2020.
- 448 [43] L. Zintgraf, K. Shiarlis, M. Igl, S. Schulze, Y. Gal, K. Hofmann, and S. Whiteson.
 449 Varibad: A very good method for bayes-adaptive deep rl via meta-learning. *arXiv preprint*
 450 *arXiv:1910.08348*, 2019.
- 451 [44] L. M. Zintgraf, L. Feng, C. Lu, M. Igl, K. Hartikainen, K. Hofmann, and S. Whiteson. Ex-
 452 ploration in approximate hyper-state space for meta reinforcement learning. In *International*
 453 *Conference on Machine Learning*, pages 12991–13001. PMLR, 2021.

454 Checklist

- 455 1. For all authors...
- 456 (a) Do the main claims made in the abstract and introduction accurately reflect the paper’s
 457 contributions and scope? [Yes]
- 458 (b) Did you describe the limitations of your work? [Yes] See Section 7
- 459 (c) Did you discuss any potential negative societal impacts of your work? [N/A] Our work
 460 is done in simulation and won’t have any negative societal impact.
- 461 (d) Have you read the ethics review guidelines and ensured that your paper conforms to
 462 them? [Yes] This work does not actually use human subjects, and is done in simulation.
 463 We have reviewed ethics guidelines and ensured that our paper conforms to them.
- 464 2. If you are including theoretical results...
- 465 (a) Did you state the full set of assumptions of all theoretical results? [N/A] Math is used
 466 as a theory/formalism, but we don’t make any provable claims about it.
- 467 (b) Did you include complete proofs of all theoretical results? [N/A]
- 468 3. If you ran experiments...
- 469 (a) Did you include the code, data, and instructions needed to reproduce the main ex-
 470 perimental results (either in the supplemental material or as a URL)? [Yes] We have
 471 included the code along with a README in the supplemental material
- 472 (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they
 473 were chosen)? [Yes] See Appendix I

- 474 (c) Did you report error bars (e.g., with respect to the random seed after running experi-
475 ments multiple times)? [Yes] All plots were created with 4 random seeds with std error
476 bars.
- 477 (d) Did you include the total amount of compute and the type of resources used (e.g., type
478 of GPUs, internal cluster, or cloud provider)? [Yes] See Appendix I
- 479 4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets...
- 480 (a) If your work uses existing assets, did you cite the creators? [Yes] Environments we
481 used are cited in section 6. Codebase used are cited in Appendix I
- 482 (b) Did you mention the license of the assets? [N/A]
- 483 (c) Did you include any new assets either in the supplemental material or as a URL? [Yes]
484 We published the code and included all environments and assets as a part of this
- 485 (d) Did you discuss whether and how consent was obtained from people whose data you're
486 using/curating? [Yes] Environments and codebases we used are open-source.
- 487 (e) Did you discuss whether the data you are using/curating contains personally identifiable
488 information or offensive content? [N/A]
- 489 5. If you used crowdsourcing or conducted research with human subjects...
- 490 (a) Did you include the full text of instructions given to participants and screenshots, if
491 applicable? [N/A]
- 492 (b) Did you describe any potential participant risks, with links to Institutional Review
493 Board (IRB) approvals, if applicable? [N/A]
- 494 (c) Did you include the estimated hourly wage paid to participants and the total amount
495 spent on participant compensation? [N/A]