

# Active Domain Randomization

**Bhairav Mehta**  
Mila, Université de Montréal

**Manfred Diaz**  
Mila, Université de Montréal

**Florian Golemo**  
Mila, Université de Montréal, ElementAI

**Christopher J. Pal**  
Mila, Polytechnique Montréal, ElementAI, CIFAR

**Liam Paull**  
Mila, Université de Montréal, CIFAR

## Abstract:

Domain randomization is a popular technique for improving domain transfer, often used in a zero-shot setting when the target domain is unknown or cannot easily be used for training. In this work, we empirically examine the effects of domain randomization on agent generalization. Our experiments show that domain randomization may lead to suboptimal, high-variance policies, which we attribute to the uniform sampling of environment parameters. We propose Active Domain Randomization, a novel algorithm that learns a parameter sampling strategy. Our method looks for the most informative environment variations within the given randomization ranges by leveraging the discrepancies of policy rollouts in randomized and reference environment instances. We find that training more frequently on these instances leads to better overall agent generalization. Our experiments across various physics-based simulated and real-robot tasks show that this enhancement leads to more robust, consistent policies.

**Keywords:** sim2real, domain randomization, reinforcement learning

## 1 Introduction

Recent trends in Deep Reinforcement Learning (DRL) exhibit a growing interest in zero-shot domain transfer, i.e. when a policy is learned in a source domain and is then tested *without finetuning* in a previously unseen target domain. Zero-shot transfer is particularly useful when the task in the target domain is inaccessible, complex, or expensive, such as gathering rollouts from a real-world robot. An ideal agent would learn to *generalize* across domains; it would accomplish the task without exploiting irrelevant features or deficiencies in the source domain (i.e., approximate physics in simulators), which may vary dramatically after transfer.

One promising approach for zero-shot transfer has been Domain Randomization (DR) [1]. In DR, we uniformly randomize environment parameters (i.e. friction, motor torque) in predefined ranges after every training episode. By randomizing everything that might vary in the target environment, the hope is that the agent will view the target domain as just another variation. However, recent works suggest that the sample complexity grows exponentially with the number of randomization parameters, even when dealing only with transfer between simulations (i.e. in Andrychowicz et al. [2] Figure 8). In addition, when using DR *unsuccessfully*, policy transfer fails, but with no clear way to understand the underlying cause. After a failed transfer, randomization ranges are tweaked heuristically via trial-and-error. Repeating this process iteratively leads to arbitrary ranges that do (or do not) lead to policy convergence without any insight into how those settings may affect the learned behavior.

In this work, we demonstrate that the strategy of *uniformly* sampling environment parameters from predefined ranges is suboptimal and propose an alternative sampling method, **Active Domain Randomization**. Active Domain Randomization (ADR), shown graphically in Figure 1, formulates

---

Correspondence to mehtabha@mila.quebec

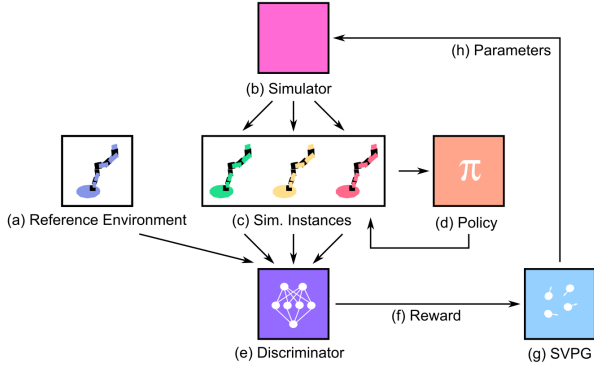


Figure 1: ADR proposes **randomized environments** (c) or simulation instances from a **simulator** (b) and rolls out an **agent policy** (d) in those instances. The **discriminator** (e) learns a **reward** (f) as a proxy for environment difficulty by distinguishing between rollouts in the **reference environment** (a) and randomized instances, which is used to train **SVPG particles** (g). The particles propose a diverse set of environments, trying to find the environment **parameters** (h) that are currently causing the agent the most difficulty.

DR as a search for randomized environments that maximize utility for the agent policy. Concretely, we aim to find environments that *currently* cause difficulties for the agent policy, dedicating more training time to these troublesome parameter settings. We cast this active search as a Reinforcement Learning (RL) problem where the ADR sampling policy is parameterized with Stein Variational Policy Gradient (SVPG) [3]. ADR focuses on problematic regions of the randomization space by learning a discriminative reward computed from discrepancies in policy rollouts generated in randomized and reference environments.

We first showcase ADR on a simple environment where the benefits of training on more challenging variations are apparent and interpretable (Figure 2). In this case, we demonstrate that ADR learns to preferentially select parameters from these more challenging parameter regions while still adapting to the policy’s current deficiencies. We then apply ADR to more complex environments and real robot settings (Figure 3) and show that even with high-dimensional search spaces and unmodeled dynamics, policies trained with ADR exhibit superior generalization and lower overall variance than their Uniform Domain Randomization (UDR) counterparts.

## 2 Preliminaries

### 2.1 Reinforcement Learning

We consider a RL framework [4] where some task  $T$  is defined by a Markov Decision Process (MDP) consisting of a state space  $S$ , action space  $A$ , state transition function  $P : S \times A \mapsto S$ , reward function  $R : S \times A \mapsto \mathbb{R}$ , and discount factor  $\gamma \in (0, 1)$ . The goal for an agent trying to solve  $T$  is to learn a policy  $\pi$  with parameters  $\theta$  that maximizes the expected total discounted reward. We define a rollout  $\tau = (s_0, a_0, \dots, s_T, a_T)$  to be the sequence of states  $s_t$  and actions  $a_t \sim \pi(a_t | s_t)$  executed by a policy  $\pi$  in the environment.

### 2.2 Stein Variational Policy Gradient

Recently, Liu et al. [3] proposed SVPG, which learns an ensemble of policies  $\mu_\phi$  in a maximum-entropy RL framework [5].

$$\max_{\mu} \mathbb{E}_{\mu} [J(\mu)] + \alpha \mathcal{H}(\mu) \quad (1)$$

with entropy  $\mathcal{H}$  being controlled by temperature parameter  $\alpha$ . SVPG uses Stein Variational Gradient Descent [6] to iteratively update an ensemble of  $N$  policies or *particles*  $\mu_\phi = \{\mu_{\phi_i}\}_{i=1}^N$  using:

$$\mu_{\phi_i} \leftarrow \mu_{\phi_i} + \frac{\epsilon}{N} \sum_{j=1}^N [\nabla_{\mu_{\phi_j}} J(\mu_{\phi_j}) k(\mu_{\phi_j}, \mu_{\phi_i}) + \alpha \nabla_{\mu_{\phi_j}} k(\mu_{\phi_j}, \mu_{\phi_i})] \quad (2)$$

with step size  $\epsilon$  and positive definite kernel  $k$ . This update rule balances exploitation (first term moves particles towards high-reward regions) and exploration (second term repulses similar policies).

### 2.3 Domain Randomization

Domain randomization (DR) is a technique to increase the generalization capability of policies trained in simulation. DR requires a prescribed set of  $N_{rand}$  simulation parameters to randomize, as well

as corresponding ranges to sample them from. A set of parameters is sampled from *randomization space*  $\Xi \subset \mathbb{R}^{N_{rand}}$ , where each randomization parameter  $\xi^{(i)}$  is bounded on a closed interval  $\{[\xi_{low}^{(i)}, \xi_{high}^{(i)}]\}_{i=1}^{N_{rand}}$ .

When a configuration  $\xi \in \Xi$  is passed to a non-differentiable simulator  $S$ , it generates an environment  $E$ . At the start of each episode, the parameters are uniformly sampled from the ranges, and the environment generated from those values is used to train the agent policy  $\pi$ .

DR may perturb any to all elements of the task  $T$ 's underlying MDP<sup>1</sup>, with the exception of keeping  $R$  and  $\gamma$  constant. DR therefore generates a set of MDPs that are superficially similar, but can vary greatly in difficulty depending on the character of the randomization. Upon transfer to the target domain, the expectation is that the agent policy has learned to generalize across MDPs, and sees the final domain as just another variation of parameters.

The most common instantiation of DR, UDR is summarized in Algorithm 2 in Appendix E. UDR generates randomized environment instances  $E_i$  by uniformly sampling  $\Xi$ . The agent policy  $\pi$  is then trained on rollouts  $\tau_i$  produced in randomized environments  $E_i$ .

### 3 Method

To start, we would like to answer the following question:

*Are all MDPs generated by uniform randomization equally useful for training?*

We consider the LunarLander-v2 environment, where the agent's task is to ground a lander in a designated zone and reward is based on the quality of landing (fuel used, impact velocity, etc). LunarLander-v2 has one main axis of randomization that we vary: the main engine strength (MES).

We aim to determine if certain environment instances (different values of the MES) are more *informative* - more efficient than others - in terms of aiding generalization. We set the total range of variation for the MES to be  $[8, 20]^2$  and find through empirical tests that lower engine strengths generate harder MDPs to solve. Under this assumption, we show the effects of *focused DR* by editing the range that the MES parameter is uniformly sampled from.

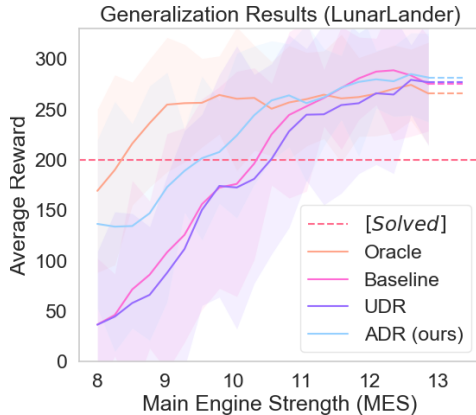


Figure 2: Agent generalization, expressed as performance across different engine strength settings in **LunarLander**. We compare the following approaches: Baseline (default environment dynamics); Uniform Domain Randomization (UDR); Active Domain Randomization (ADR, our approach); and Oracle (a handpicked randomization range of MES  $[8, 11]$ ). ADR achieves for near-expert levels of generalization, while both Baseline and UDR fail to solve lower MES tasks.

beforehand and (2) DR is used mostly when the space of randomized parameters is high-dimensional

<sup>1</sup>The effects of DR on action space  $A$  are usually implicit or are carried out on the simulation side.

<sup>2</sup>Default MES is 13;  $MES \leq 7.5$  is unsolvable when all other parameters remain constant.

We train multiple agents on different randomization ranges for MES, which define what types of environments the agent is exposed to during training. Figure 2 shows the final generalization performance of each agent by sampling randomly from the entire randomization range of  $[8, 20]$  and rolling out the policy in the generated environments. We see that, in this case, focusing on harder MDPs improves generalization as compared to uniformly sampling the whole space, even when the evaluation environment is outside of the training distribution.

#### 3.1 Problem Formulation

The experiment in the previous section shows that preferential training on more informative environments provides tangible benefits in terms of agent generalization. However, in general, finding these informative environments is difficult because: (1) It is rare that such intuitively *hard* MDP instances or parameter ranges are known

or noninterpretable. As a result, we propose an algorithm for finding environment instances that maximize *utility*, or provide the most improvement (in terms of generalization) to our agent policy when used for training.

### 3.2 Active Domain Randomization

Drawing analogies with Bayesian Optimization (BO) literature, one can consider the randomization space as a search space. Traditionally, in BO, the search for where to evaluate an objective is informed by acquisition functions, which trade off exploitation of the objective with exploration in the uncertain regions of the space [7]. However, unlike the stationary objectives seen in BO, training the agent policy renders our optimization non-stationary: the environment with highest utility at time  $t$  is likely not the same as the maximum utility environment at time  $t + 1$ . This requires us to redefine the notion of an acquisition function while simultaneously dealing with BO’s deficiencies with higher-dimensional inputs [8].

---

#### Algorithm 1 Active Domain Randomization

---

```

1: Input:  $\Xi$ : Randomization space,  $S$ : Simulator,  $\xi_{ref}$ :
   reference parameters
2: Initialize  $\pi_\theta$ : agent policy,  $\mu_\phi$ : SVPG particles,  $D_\psi$ :
   discriminator,  $E_{ref} \leftarrow S(\xi_{ref})$ : reference environment
3: while not max.timesteps do
4:   for each particle  $\mu_\phi$  do
5:     rollout  $\xi_i \sim \mu_\phi(\cdot)$ 
6:   end for
7:   for each  $\xi_i$  do
8:     // Generate, rollout in randomized env.
9:      $E_i \leftarrow S(\xi_i)$ 
10:    rollout  $\tau_i \sim \pi_\theta(\cdot; E_i)$ ,  $\tau_{ref} \sim \pi_\theta(\cdot; E_{ref})$ 
11:     $\mathcal{T}_{rand} \leftarrow \mathcal{T}_{rand} \cup \tau_i$ ;  $\mathcal{T}_{ref} \leftarrow \mathcal{T}_{ref} \cup \tau_{ref}$ 
12:  end for
13:  // Calculate reward for each proposed environment
14:  for each  $\tau_i \in \mathcal{T}_{rand}$  do
15:    Calculate  $r_i$  for  $\xi_i / E_i$  (Eq. (3))
16:  end for
17:  // Gradient Updates
18:  with  $\mathcal{T}_{rand}$  update:
19:     $\theta \leftarrow \theta + \nu \nabla_\theta J(\pi_\theta)$ 
20:  Update particles using Eq. (2)
21:  Update  $D_\psi$  with  $\tau_i$  and  $\tau_{ref}$  using SGD.
22: end while

```

---

To this end, we propose ADR, summarized in Algorithm 1 and Figure 1<sup>3</sup>. ADR provides a framework for manipulating a more general analog of an *acquisition function*, selecting the most informative MDPs for the agent within the randomization space. By formulating the search as an RL problem, ADR learns a policy  $\mu_\phi$  where states are proposed randomization configurations  $\xi \in \Xi$  and actions are continuous changes to those parameters.

We learn a discriminator-based reward for  $\mu_\phi$ , similar the reward seen in Eysenbach et al. [9]:

$$r_D = \log D_\psi(y|\tau_i \sim \pi(\cdot; E_i)) \quad (3)$$

where  $y$  is a boolean variable denoting the discriminator’s prediction of which type of environment (a randomized environment  $E_i$  or reference environment  $E_{ref}$ ) the trajectory  $\tau_i$  was

generated from. We assume that the  $E_{ref} = S(\xi_{ref})$  is provided with the original task definition.

Intuitively, we reward the policy  $\mu_\phi$  for finding regions of the randomization space that produce environment instances where the *same* agent policy  $\pi$  acts differently than in the reference environment. The agent policy  $\pi$  sees and trains *only* on the randomized environments (as it would in traditional DR), using the environment’s task-specific reward for updates. As the agent improves on the proposed, problematic environments, it becomes more difficult to differentiate whether any given state transition was generated from the reference or randomized environment. Thus, ADR can find what parts of the randomization space the agent is currently performing poorly on, and can *actively* update its sampling strategy throughout the training process.

## 4 Results

### 4.1 Experiment Details

To test ADR, we experiment on OpenAI Gym environments [10] across various tasks, both simulated and real: (a) LunarLander-v2, a 2 degrees of freedom (DoF) environment where the agent has to

---

<sup>3</sup>We provide a detailed walkthrough of the algorithm in Appendix A.

softly land a spacecraft, implemented in Box2D (detailed in Section 3.2), **(b)** Pusher-3Dof-v0, a 3 DoF arm from Haarnoja et al. [11] that has to push a puck to a target, implemented in Mujoco [12], and **(c)** ErgoReacher-v0, a 4 DoF arm from Golemo et al. [13] which has to touch a goal with its end effector, implemented in the Bullet Physics Engine [14]. For sim2real experiments, we recreate this environment setup on a real Poppy Ergo Jr. robot [15] shown in Figure 3 (a) and (b), and also create **(d)** ErgoPusher-v0 an environment similar to Pusher-3Dof-v0 with a real robot analog seen in Figure 3 (c) and (d). We provide a detailed account of the randomized parameters in each environment in Table 1 in Appendix F.

All simulated experiments are run with five seeds each with five random resets, **totaling 25 independent trials per evaluation point**. All experimental results are plotted mean-averaged with one standard deviation shown. Detailed experiment information can be found in Appendix H.

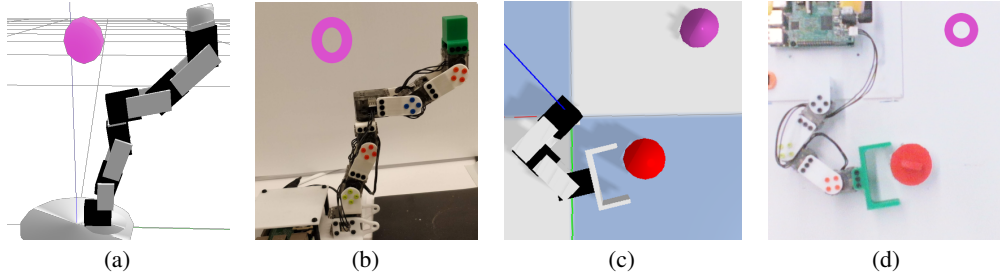


Figure 3: Along with simulated environments, we display ADR on zero-shot transfer tasks onto real robots.

## 4.2 Toy Experiments

To investigate whether ADR’s learned sampling strategy provides a tangible benefit for agent generalization, we start by comparing it against traditional DR (labeled as UDR) on LunarLander-v2 and vary only the main engine strength (MES). In Figure 2, we see that ADR approaches expert-levels of generalization whereas UDR fails to generalize on lower MES ranges.

We compare the learning progress for the different methods on the *hard* environment instances ( $\xi_{MES} \sim U[8, 11]$ ) in Figure 4(a). ADR significantly outperforms both the baseline (trained only on MES of 13) and the UDR agent (trained seeing environments with  $\xi_{MES} \sim U[8, 20]$ ) in terms of performance.

Figures 4(b) and 4(c) showcase the adaptability of ADR by showing generalization and sampling distributions at various stages of training. ADR samples approximately uniformly for the first 650K steps, but then finds a deficiency in the policy on higher ranges of the MES. As those areas become more frequently sampled between 650K-800K steps, the agent learns to solve all of the higher-MES environments, as shown by the generalization curve for 800K steps. As a result, the discriminator is no longer able to differentiate reference and randomized trajectories from the higher MES regions, and starts to reward environment instances generated in the lower end of the MES range, which improves generalization towards the completion of training.

## 4.3 Randomization in High Dimensions

If the intuitions that drive ADR are correct, we should see increased benefit of a learned sampling strategy with larger  $N_{rand}$  due to the increasing sparsity of *informative* environments when sampling uniformly. We first explore ADR’s performance on Pusher3Dof-v0, an environment where  $N_{rand} = 2$ . Both randomization dimensions (puck damping, puck friction loss) affect whether or not the puck retains momentum and continues to slide after making contact with the agent’s end effector. Lowering the values of these parameters *simultaneously* creates an intuitively-harder environment, where the puck continues to slide after being hit. In the reference environment, the puck retains no momentum and must be continuously pushed in order to move. We qualitatively visualize the effect of these parameters on puck sliding in Figure 5(a).

From Figure 5(b), we see ADR’s improved robustness to *extrapolation* - or when the target domain lies *outside the training region*. We train two agents, one using ADR and one using UDR, and show them only the training regions encapsulated by the dark, outlined box in the top-right of Figure 5(a).

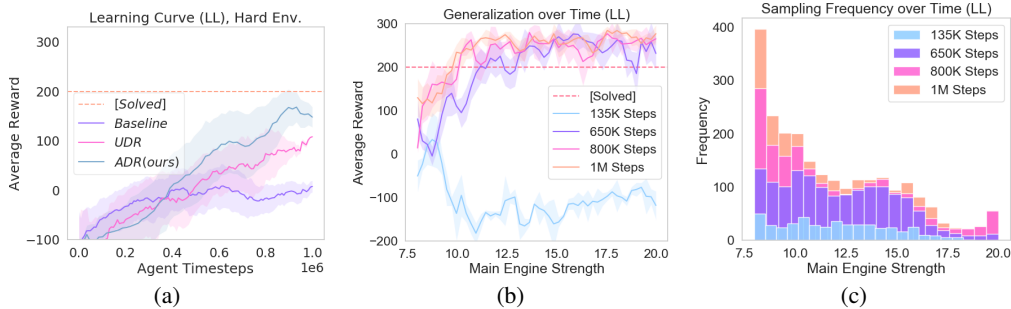


Figure 4: Learning curves over time in **LunarLander**. Higher is better. (a) Performance on particularly difficult settings - our approach outperforms both the policy trained on a single simulator instance (“baseline”) and the UDR approach. (b) Agent generalization in **LunarLander** over time during training when using ADR. (c) Adaptive sampling visualized. ADR, seen in (b) and (c), adapts to where the agent is struggling the most, improving generalization performance by end of training.

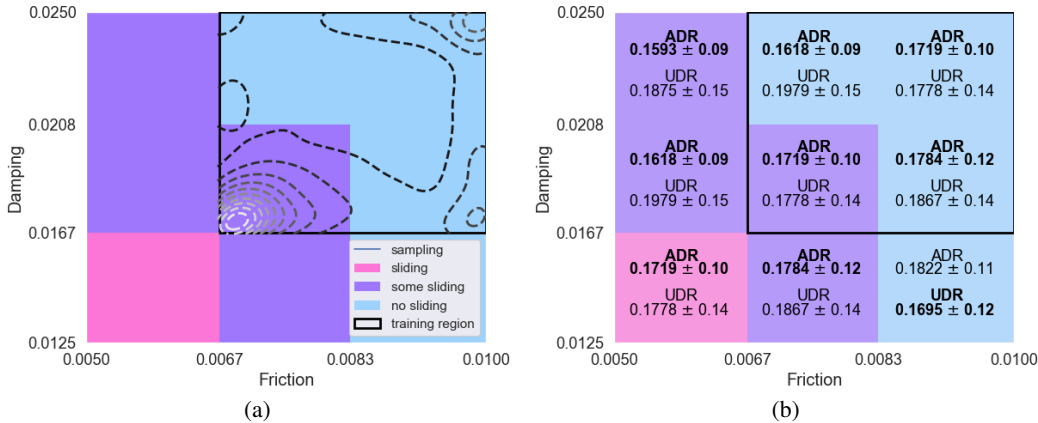


Figure 5: In **Pusher-3Dof**, the environment dynamics are characterized by friction and damping of the sliding puck, where sliding correlates with the difficulty of the task (as highlighted by cyan, purple, and pink - from easy to hard). (a) During training, the algorithm only had access to a limited, easier range of dynamics (black outlined box in the upper right). (b) Performance measured by distance to target, lower is better.

Qualitatively, only 25% of the environments have dynamics which cause the puck to slide, which are the hardest environments to solve in the training region. We see that from the sampling histogram overlaid on Figure 5(a) that ADR prioritizes the single, harder purple region more than the light blue regions, allowing for better generalization to the unseen test domains, as shown in Figure 5(b). ADR outperforms UDR in all but one test region and produces policies with less variance than their UDR counterparts.

#### 4.4 Randomization in *Uninterpretable* Dimensions

We further show the significance of ADR over UDR on **ErgoReacher-v0**, where  $N_{rand} = 8$ . It is now impossible to infer intuitively which environments are *hard* due to the complex interactions between the eight randomization parameters (gains and maximum torques for each joint). For demonstration purposes, we test extrapolation by creating a held-out target environment with extremely low values for torque and gain, which causes certain states in the environment to lead to *catastrophic* failure - gravity pulls the robot end effector down, and the robot is not strong enough to pull itself back up. We show an example of an agent getting trapped in a catastrophic failure state in Figure 11, Appendix F.1.

To generalize effectively, the sampling policy should prioritize environments with lower torque and gain values in order for the agent to operate in such states precisely. However, since the hard evaluation environment is not seen during training, ADR must learn to prioritize the hardest environments that it can see, while still learning behaviors that can operate well across the entire training region.



From Figure 6(a) (learning curves for Pusher3Dof-v0 on the unseen, hard environment - the pink square in Figure 5) and 6(b) (learning curves for ErgoReacher-v0 on unseen, hard environment), we observe the detrimental effects of uniform sampling. In Pusher3Dof-v0, we see that UDR *unlearns* the good behaviors it acquired in the beginning of training. When training neural networks in both supervised and reinforcement learning settings, this phenomenon has been dubbed as *catastrophic forgetting* [16]. ADR seems to exhibit this slightly (leading to "hills" in the curve), but due to the adaptive nature the algorithm, it is able to adjust quickly and retain better performance across all environments.

UDR’s high variance on ErgoReacher-v0 highlights another issue: by continuously training on a random mix of hard and easy MDP instances, both beneficial and detrimental agent behaviors can be learned and unlearned throughout training. This mixing can lead to high-variance, inconsistent, and unpredictable behavior upon transfer. By focusing on those harder environments and allowing the definition of *hard* to adapt over time, ADR shows more consistent performance and better overall generalization than UDR in all environments tested.

#### 4.5 Sim2Real Transfer Experiments

In *sim2real* (simulation to reality) transfer, many policies fail due to unmodeled dynamics within the simulators, as policies may have overfit to or exploited simulation-specific details of their training environments. While the deficiencies and high variance of UDR are clear even in simulated environments, one of the most impressive results of domain randomization was zero-shot transfer out of simulation onto robots. However, we find that the same issues of unpredictable performance apply to UDR-trained policies in the real world as well.

We take each method’s (ADR and UDR) five independent simulation-trained policies on both ErgoReacher-v0 and ErgoPusher-v0 and transfer them without fine tuning onto the real robot. We rollout only the final policy on the robot, and show performance in Figure 7. To evaluate generalization, we alter the environment manually: on ErgoReacher-v0, we change the values of the torques (higher torque means the arm moves at higher speed and accelerates faster); on ErgoPusher-v0, we change the friction of the sliding puck (slippery or rough). For each environment, we evaluate each of the policies with 25 random goals (125 independent evaluations per method per environment setting).

Even in zero-shot transfer tasks onto real robots, ADR policies obtain overall better or similar performance than UDR policies trained in the same conditions. More importantly, ADR policies are more consistent and display lower spread across all environments, which is crucial when safely evaluating reinforcement learning policies on real-world robots.

## 5 Related Work

### 5.1 Dynamic and Adversarial Simulators

Simulators have played a crucial role in transferring learned policies onto real robots, and many different strategies have been proposed. Randomizing simulation parameters for better generalization or transfer performance is a well-established idea in evolutionary robotics [17, 18], but recently has emerged as an effective way to perform zero-shot transfer of deep reinforcement learning policies in difficult tasks [2, 1, 19, 20].

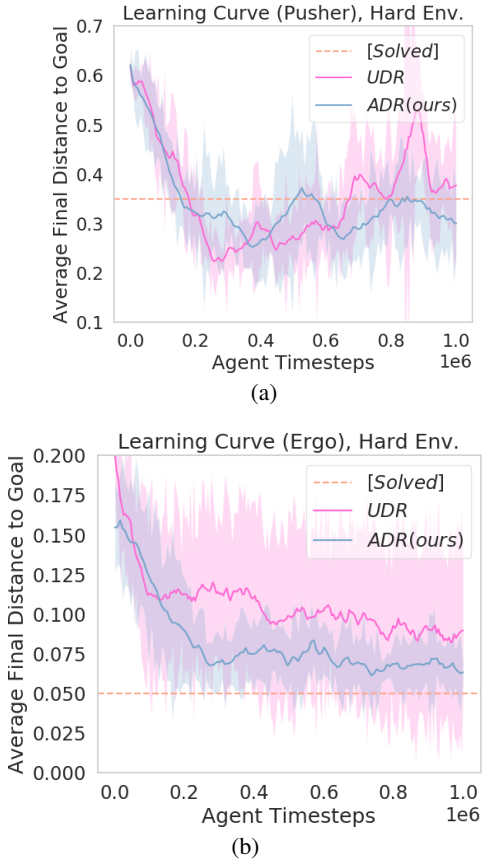
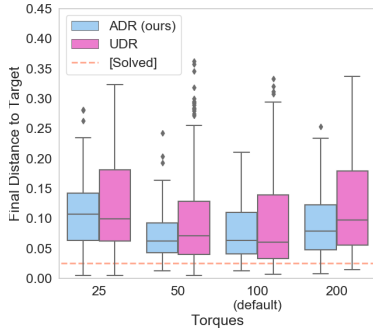
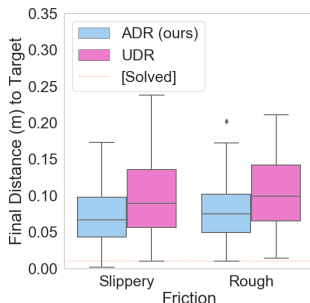


Figure 6: Learning curves over time in (a) **Pusher3Dof-v0** and (b) **ErgoReacher** on held-out, difficult environment settings. Our approach outperforms both the policy trained with the UDR approach both in terms of performance and variance.



(a)



(b)

Figure 7: Zero-shot transfer onto real robots (a) **ErgoReacher** and (b) **ErgoPusher**. In both environments, we assess generalization by manually changing torque strength and puck friction respectively.

to be labeled. The chosen samples are labelled by an oracle and sent back to the model for use. Similarly, ADR searches for what environments may be most useful to the agent at any given time. Active learners, like BO methods discussed in Section 3, often require an acquisition function (derived from a notion of model uncertainty) to chose the next sample. Since ADR handles this decision through the explore-exploit framework of RL and the  $\alpha$  in SVPG, ADR sidesteps the well-known scalability issues of both active learning and BO [26].

Recently, Toneva et al. [27] showed that certain examples in popular computer vision datasets are harder to learn, and that some examples are forgotten much quicker than others. We explore the same phenomenon in the space of MDPs defined by our randomization ranges, and try to find the "examples" that cause our agent the most trouble. Unlike in active learning or Toneva et al. [27], we have no oracle or supervisory loss signal in RL, and instead attempt to learn a proxy signal for ADR via a discriminator.

## 6 Conclusion

In this work, we highlight failure cases of traditional domain randomization, and propose active domain randomization (ADR), a general method capable of finding the most informative parts of the randomization parameter space for a reinforcement learning agent to train on. ADR does this by posing the search as a reinforcement learning problem, and optimizes for the most informative environments using a learned reward and multiple policies. We show on a wide variety of simulated environments that this method efficiently trains agents with better generalization than traditional domain randomization, extends well to high dimensional parameter spaces, and produces more robust policies when transferring to the real world.

Learnable simulations are also an effective way to adapt a simulation to a particular target environment. Chebotar et al. [21] and Ruiz et al. [22] use RL for effective transfer by learning parameters of a simulation that accurately describes the target domain, but require the target domain for reward calculation, which can lead to overfitting. In contrast, our approach requires no target domain, but rather only a reference domain (the default simulation parameters) and a general range for each parameter. ADR encourages diversity, and as a result gives the agent a wider variety of experience. In addition, unlike Chebotar et al. [21], our method does not requires carefully-tuned (co-)variances or task-specific cost functions. Concurrently, Khirodkar et al. [23] also showed the advantages of learning adversarial simulations and disadvantages of purely uniform randomization distributions in object detection tasks.

To improve policy robustness, Robust Adversarial Reinforcement Learning (RARL) Pinto et al. [24] jointly trains both an agent and an adversary who applies environment forces that disrupt the agent’s task progress. ADR removes the zero-sum game dynamics, which have been known to decrease training stability [25]. More importantly, our method’s final outputs - the SVPG-based sampling strategy and discriminator - are reusable and can be used to train new agents as shown in Appendix D, whereas a trained RARL adversary would overpower any new agent and impede learning progress.

## 5.2 Active Learning and Informative Samples

Active learning methods in supervised learning try to construct a *representative*, sometimes time-variant, dataset from a large pool of unlabelled data by proposing elements



## Acknowledgements

The authors gratefully acknowledge the Natural Sciences and Engineering Research Council of Canada (NSERC), the Fonds de Recherche Nature et Technologies Quebec (FQRNT) and the Open Philanthropy Project for supporting this work. In addition, the authors would like to thank Kyle Kastner and members of the REAL Lab for their helpful comments, as well as Nvidia for donating a DGX-1 used for this research. BM would like to thank IVADO.

## References

- [1] J. Tobin, R. Fong, A. Ray, J. Schneider, W. Zaremba, and P. Abbeel. Domain randomization for transferring deep neural networks from simulation to the real world. In *Intelligent Robots and Systems (IROS), 2017 IEEE/RSJ International Conference on*. IEEE, 2017.
- [2] M. Andrychowicz, B. Baker, M. Chociej, R. Jozefowicz, B. McGrew, J. Pachocki, A. Petron, M. Plappert, G. Powell, A. Ray, et al. Learning dexterous in-hand manipulation. *arXiv preprint arXiv:1808.00177*, 2018.
- [3] Y. Liu, P. Ramachandran, Q. Liu, and J. Peng. Stein variational policy gradient, 2017.
- [4] R. S. Sutton and A. G. Barto. *Reinforcement Learning: An introduction*. MIT Press, 2018.
- [5] B. D. Ziebart. *Modeling Purposeful Adaptive Behavior with the Principle of Maximum Causal Entropy*. PhD thesis, CMU, 2010.
- [6] Q. Liu and D. Wang. Stein variational gradient descent: A general purpose bayesian inference algorithm. In D. D. Lee, M. Sugiyama, U. V. Luxburg, I. Guyon, and R. Garnett, editors, *Advances in Neural Information Processing Systems 29*. 2016.
- [7] E. Brochu, V. M. Cora, and N. de Freitas. A tutorial on bayesian optimization of expensive cost functions, with application to active user modeling and hierarchical reinforcement learning. *CoRR*, abs/1012.2599, 2010. URL <http://arxiv.org/abs/1012.2599>.
- [8] Z. Wang, M. Zoghi, F. Hutter, D. Matheson, and N. De Freitas. Bayesian optimization in high dimensions via random embeddings. In *Proceedings of the Twenty-Third International Joint Conference on Artificial Intelligence, IJCAI '13*, 2013.
- [9] B. Eysenbach, A. Gupta, J. Ibarz, and S. Levine. Diversity is all you need: Learning skills without a reward function. *arXiv preprint arXiv:1802.06070*, 2018.
- [10] G. Brockman, V. Cheung, L. Pettersson, J. Schneider, J. Schulman, J. Tang, and W. Zaremba. Openai gym, 2016.
- [11] T. Haarnoja, V. Pong, A. Zhou, M. Dalal, P. Abbeel, and S. Levine. Composable deep reinforcement learning for robotic manipulation. *arXiv preprint arXiv:1803.06773*, 2018.
- [12] E. Todorov, T. Erez, and Y. Tassa. Mujoco: A physics engine for model-based control. In *IROS*. IEEE, 2012.
- [13] F. Golemo, A. A. Taiga, A. Courville, and P.-Y. Oudeyer. Sim-to-real transfer with neural-augmented robot simulation. In *Conference on Robot Learning*, 2018.
- [14] E. Coumans. Bullet physics simulation. In *ACM SIGGRAPH 2015 Courses*, SIGGRAPH '15, New York, NY, USA, 2015. ACM.
- [15] M. Lapeyre. *Poppy: open-source, 3D printed and fully-modular robotic platform for science, art and education*. Theses, Université de Bordeaux, Nov. 2014. URL <https://hal.inria.fr/tel-01104641>.
- [16] J. Kirkpatrick, R. Pascanu, N. C. Rabinowitz, J. Veness, G. Desjardins, A. A. Rusu, K. Milan, J. Quan, T. Ramalho, A. Grabska-Barwinska, D. Hassabis, C. Clopath, D. Kumaran, and R. Hadsell. Overcoming catastrophic forgetting in neural networks. *CoRR*, abs/1612.00796, 2016. URL <http://arxiv.org/abs/1612.00796>.
- [17] J. C. Zagal, J. Ruiz-del Solar, and P. Vallejos. Back to reality: Crossing the reality gap in evolutionary robotics. *IFAC Proceedings Volumes*, 37(8), 2004.

- [18] J. Bongard and H. Lipson. Once more unto the breach: Co-evolving a robot and its simulator. In *Proceedings of the Ninth International Conference on the Simulation and Synthesis of Living Systems (ALIFE9)*, 2004.
- [19] X. B. Peng, M. Andrychowicz, W. Zaremba, and P. Abbeel. Sim-to-real transfer of robotic control with dynamics randomization. *2018 IEEE International Conference on Robotics and Automation (ICRA)*, May 2018.
- [20] F. Sadeghi and S. Levine. (cad)\$^2\$rl: Real single-image flight without a single real image. *CoRR*, abs/1611.04201, 2016. URL <http://arxiv.org/abs/1611.04201>.
- [21] Y. Chebotar, A. Handa, V. Makoviychuk, M. Macklin, J. Issac, N. Ratliff, and D. Fox. Closing the sim-to-real loop: Adapting simulation randomization with real world experience. *arXiv preprint arXiv:1810.05687*, 2018.
- [22] N. Ruiz, S. Schulter, and M. Chandraker. Learning to simulate, 2018.
- [23] R. Khirodkar, D. Yoo, and K. M. Kitani. Vadra: Visual adversarial domain randomization and augmentation. *arXiv preprint arXiv:1812.00491*, 2018.
- [24] L. Pinto, J. Davidson, R. Sukthankar, and A. Gupta. Robust adversarial reinforcement learning, 2017.
- [25] L. Mescheder, A. Geiger, and S. Nowozin. Which training methods for GANs do actually converge? In *International Conference on Machine Learning*, 2018.
- [26] S. Tong. *Active Learning: Theory and Applications*. PhD thesis, 2001. AAI3028187.
- [27] M. Toneva, A. Sordoni, R. T. d. Combes, A. Trischler, Y. Bengio, and G. J. Gordon. An empirical study of example forgetting during deep neural network learning. *arXiv preprint arXiv:1812.05159*, 2018.
- [28] T. P. Lillicrap, J. J. Hunt, A. Pritzel, N. Heess, T. Erez, Y. Tassa, D. Silver, and D. Wierstra. Continuous control with deep reinforcement learning. *arXiv preprint arXiv:1509.02971*, 2015.
- [29] D. P. Kingma and J. Ba. Adam: A method for stochastic optimization, 2014.
- [30] S. Fujimoto, H. van Hoof, and D. Meger. Addressing function approximation error in actor-critic methods. In *International Conference on Machine Learning*, 2018.
- [31] V. Mnih, A. P. Badia, M. Mirza, A. Graves, T. Lillicrap, T. Harley, D. Silver, and K. Kavukcuoglu. Asynchronous methods for deep reinforcement learning. In *International Conference on Machine Learning*, pages 1928–1937, 2016.
- [32] T. Gangwani, Q. Liu, and J. Peng. Learning self-imitating diverse policies. In *International Conference on Learning Representations*, 2019.

## Appendix A Architecture Walkthrough

In this section, we walk through the diagram shown in Figure 1. All line references refer to Algorithm 1.

### A.0.1 SVPG Sampler

To encourage sufficient exploration in high dimensional randomization spaces, we parameterize  $\mu_\phi$  with SVPG. Since each particle proposes its own environment settings  $\xi_i$  (lines 4-6, Figure 1h), all of which are passed to the agent for training, the agent policy benefits from the same environment variety seen in UDR. However, unlike UDR,  $\mu_\phi$  can use the learned reward to focus on problematic MDP instances while still being efficiently parallelizable.

### A.0.2 Simulator

After receiving each particle’s proposed parameter settings  $\xi_i$ , we generate randomized environments  $E_i = S(\xi_i)$  (line 9, Figure 1b).

### A.0.3 Generating Trajectories

We proceed to train the agent policy  $\pi$  on the randomized instances  $E_i$ , just as in UDR. We roll out  $\pi$  on each randomized instance  $E_i$  and store each trajectory  $\tau_i$ . For every randomized trajectory generated, we use the *same* policy to collect and store a reference trajectory  $\tau_{ref}$  by rolling out  $\pi$  in the default environment  $E_{ref}$  (lines 10-12, Figure 1a, c). We store all trajectories (lines 11-12) as we will use them to score each parameter setting  $\xi_i$  and update the discriminator.

The agent policy is a *black box*: although in our experiments we train  $\pi$  with Deep Deterministic Policy Gradients [28], the policy can be trained with any other on or off-policy algorithm by introducing only minor changes to Algorithm 1 (lines 13-17, Figure 1d).

### A.0.4 Scoring Environments

We now generate a score for each environment (lines 20-22) using each stored randomized trajectory  $\tau_i$  by passing them through the discriminator  $D_\psi$ , which predicts the type of environment (reference or randomized) each trajectory was generated from. We use this score as a reward to update each SVPG particle using Equation 2 (lines 24-26, Figure 1f).

After scoring each  $\xi_i$  according to Equation 3, we use the randomized and reference trajectories to train the discriminator (lines 28-30, Figure 1e).

## Appendix B Learning Curves for Reference Environments

For space concerns, we show only the *hard* generalization curves for all environments in the main document. For completeness, we include learning curves on the reference environment here.

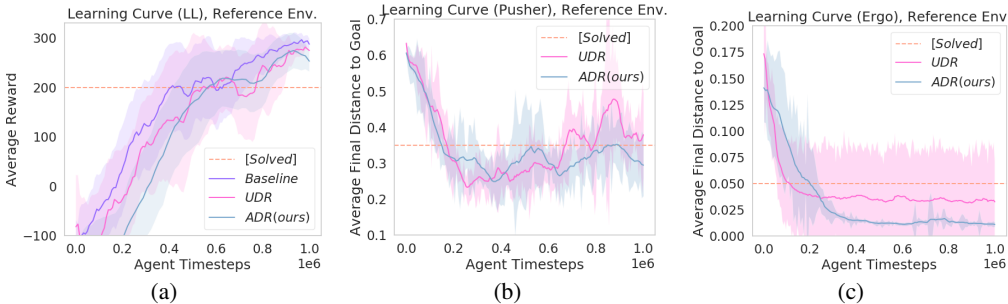


Figure 8: Learning curves over time reference environments. (a) LunarLander (b) Pusher-3Dof (c) ErgoReacher.

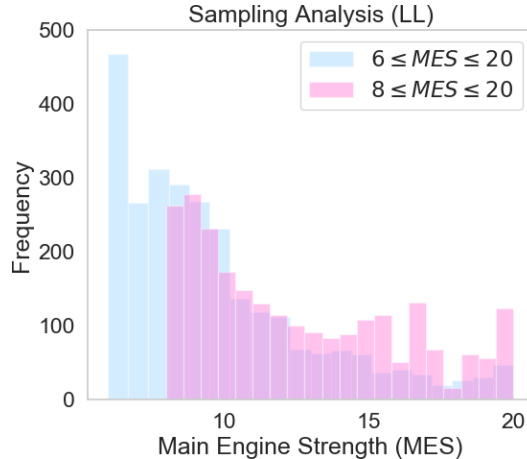


Figure 9: Sampling frequency across engine strengths when varying the randomization ranges. The updated, red distribution shows a much milder unevenness in the distribution, while still learning to focus on the harder instances.

### Appendix C Interpretability Benefits of ADR

One of the secondary benefits of ADR is its insight into incompatibilities between the task and randomization ranges. We demonstrate the simple effects of this phenomenon in a one-dimensional LunarLander-v2, where we only randomize the main engine strength. Our initial experiments varied this parameter between 6 and 20, which lead to ADR learning degenerate agent policies by learning to propose the lopsided blue distribution in Figure 9. Upon inspection of the simulation, we see that when the parameter has a value of less than approximately 8, the task becomes almost impossible to solve due to the other environment factors (in this case the lander always hits the ground too fast, which it is penalized for).

After adjusting the parameter ranges to more sensible values, we see a better sampled distribution in pink, which still gives more preference to the hard environments in the lower engine strength range. Most importantly, ADR allows for analysis that is both *focused* - we know exactly what part of the simulation is causing trouble - and *pre-transfer*, i.e. done before a more expensive experiment such as real robot transfer has taken place. With UDR, the agents would be equally trained on these degenerate environments, leading to policies with potentially undefined behavior (or, as seen in Section 4.4, unlearn good behaviors) in these truly out-of-distribution simulations.

### Appendix D Bootstrapping Training of New Agents

Unlike DR, ADR’s learned sampling strategy and discriminator can be reused to train new agents from scratch. To test the transferability of the sampling strategy, we first train an instance of ADR on LunarLander-v2, and then extract the SVPG particles and discriminator. We then replace the agent policy with a random network initialization, and once again train according to the details in Section 4.1. From Figure 10(a), it can be seen that the bootstrapped agent generalization is even better than the one learned with ADR from scratch. However, its training speed on the default environment ( $\xi_{MES} = 13$ ) is relatively lower.

### Appendix E Uniform Domain Randomization

Here we review the algorithm for Uniform Domain Randomization (UDR), first proposed in [1], shown in Algorithm 2.

### Appendix F Environment Details

Please see Table 1.

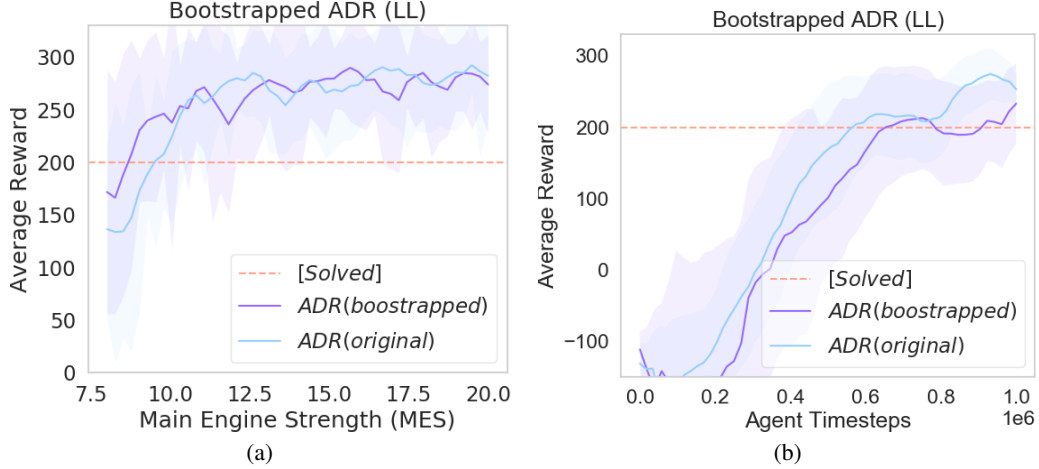


Figure 10: Generalization and default environment learning progression on LunarLander-v2 when using ADR to bootstrap a new policy. Higher is better.

---

### Algorithm 2 Uniform Sampling Domain Randomization

---

```

1: Input:  $\Xi$ : Randomization space,  $S$ : Simulator
2: Initialize  $\pi_\theta$ : agent policy
3: for each episode do
4:   // Uniformly sample parameters
5:   for  $i = 1$  to  $N_{rand}$  do
6:      $\xi^{(i)} \sim U[\xi_{low}, \xi_{high}]$ 
7:   end for
8:   // Generate, rollout in randomized env.
9:    $E_i \leftarrow S(\xi_i)$ 
10:  rollout  $\tau_i \sim \pi_\theta(\cdot; E_i)$ 
11:   $\mathcal{T}_{rand} \leftarrow \mathcal{T}_{rand} \cup \tau_i$ 
12:  for each gradient step do
13:    // Agent policy update
14:    with  $\mathcal{T}_{rand}$  update:
15:       $\theta \leftarrow \theta + \nu \nabla_\theta J(\pi_\theta)$ 
16:  end for
17: end for

```

---

### F.1 Catastrophic Failure States in ErgoReacher

In Figure 11, we show an example progression to a *catastrophic failure state* in the held-out, simulated target environment of ErgoReacher-v0, with extremely low torque and gain values.

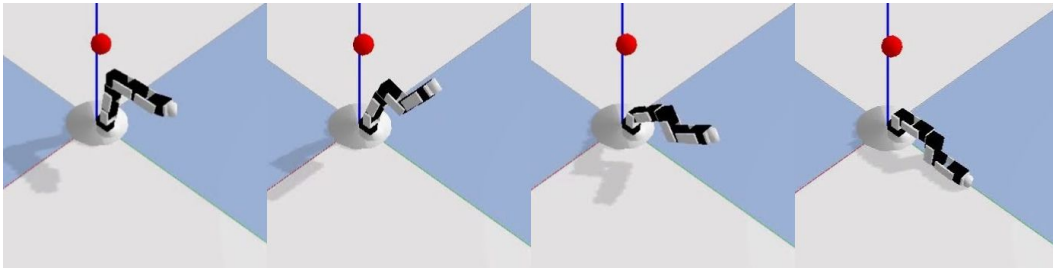


Figure 11: An example progression (left to right) of an agent moving to a catastrophic failure state (Panel 4) in the hard ErgoReacher-v0 environment.



Environment	$N_{rand}$	Types of Randomizations	Train Ranges	Test Ranges
LunarLander-v2	1	Main Engine Strength	[8, 20]	[8, 11]
Pusher-3DOF-v0	2	Puck Friction Loss & Puck Joint Damping	$[0.67, 1.0] \times \text{default}$	$[0.5, 0.67] \times \text{default}$
ErgoPusher-v0	2	Puck Friction Loss & Puck Joint Damping	$[0.67, 1.0] \times \text{default}$	$[0.5, 0.67] \times \text{default}$
ErgoReacher-v0	8	Joint Damping	$[0.3, 2.0] \times \text{default}$	$0.2 \times \text{default}$
		Joint Max Torque	$[1.0, 4.0] \times \text{default}$	default

Table 1: We summarize the environments used, as well as characteristics about the randomizations performed in each environment.

## Appendix G Untruncated Plots for Lunar Lander

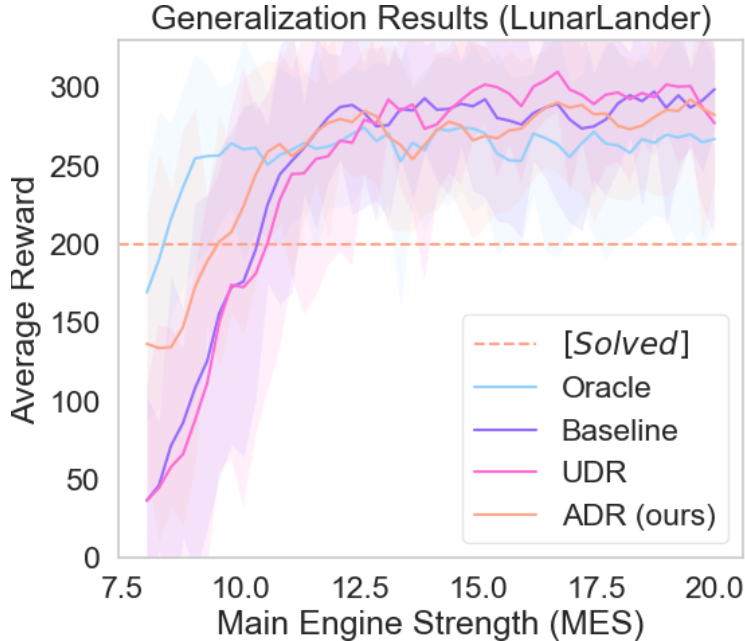


Figure 12: Generalization on LunarLander-v2 for an expert interval selection, ADR, and UDR. Higher is better.

All policies on Lunar Lander described in our paper receive a *Solved* score when the engine strengths are above 12, which is why truncated plots are shown in the main document. For clarity, we show the full, untruncated plot in Figure 12.

## Appendix H Network Architectures and Experimental Hyperparameters

All experiments can be reproduced using our Github repository<sup>4</sup>.

All of our experiments use the same network architectures and experiment hyperparameters, except for the number of particles  $N$ . For any experiment with LunarLander-v2, we use  $N = 10$ . For both other environments, we use  $N = 15$ . All other hyperparameters and network architectures remain constant, which we detail below. All networks use the Adam optimizer [29].

We run Algorithm 1 until 1 million *agent timesteps* are reached - i.e. the agent policy takes 1M steps in the randomized environments. We also cap each episode off a particular number of timesteps according to the documentation associated with [10]. In particular, LunarLander-v2 has an episode time limit of 1000 environment timesteps, whereas both Pusher-3DOF-v0 and ErgoReacher-v0 use an episode time limit of 100 timesteps.

<sup>4</sup>Code link will be updated after review.

For our agent policy, we use an implementation of DDPG (particularly, `OurDDPG.py`) from the Github repository associated with [30]. The actor and critic both have two hidden layers of 400 and 300 neurons respectively, and use ReLU activations. Our discriminator-based rewarder is a two-layer neural network, both layers having 128 neurons. The hidden layers use `tanh` activation, and the network outputs a `sigmoid` for prediction.

The agent particles in SVPG are parameterized by a two-layer actor-critic architecture, both layers in both networks having 100 neurons. We use Advantage Actor-Critic (A2C) to calculate unbiased and low variance gradient estimates. All of the hidden layers use `tanh` activation and are orthogonally initialized, with a learning rate of 0.0003 and discount factor  $\gamma = 0.99$ . They operate on a  $\mathbb{R}^{N_{rand}}$  continuous space, with each axis bounded between  $[0, 1]$ . We allow for set the max step length to be 0.05, and every 50 timesteps, we reset each particle and randomly initialize its state using a  $N_{rand}$ -dimensional uniform distribution. We use a temperature  $\alpha = 10$  with an RBF-Kernel as was done in [3]. In our work we use an Radial Basis Function (RBF) kernel with median baseline as described in Liu et al. [3] and an A2C policy gradient estimator [31], although both the kernel and estimator could be substituted with alternative methods [32]. To ensure diversity of environments throughout training, we always roll out the SVPG particles using a non-deterministic sample.

For DDPG, we use a learning rate  $\nu = 0.001$ , target update coefficient of 0.005, discount factor  $\gamma = 0.99$ , and batch size of 1000. We let the policy run for 1000 steps before any updates, and clip the max action of the actor between  $[-1, 1]$  as prescribed by each environment.

Our discriminator-based reward generator is a network with two, 128-neuron layers with a learning rate of .0002 and a binary cross entropy loss (i.e. is this a *randomized* or *reference* trajectory). To calculate the reward for a trajectory for any environment, we split each trajectory into its  $(s_t, a_t, s_{t+1})$  constituents, pass each tuple through the discriminator, and average the outputs, which is then set as the reward for the trajectory. Our batch size is set to be 128, and most importantly, as done in [9], we calculate the reward for examples before using those same examples to train the discriminator.