
Flow-based Domain Randomization for Learning and Sequencing Robotic Skills

Aidan Curtis¹ Eric Li¹ Michael Noseworthy¹ Nishad Gothoskar¹ Sachin Chitta² Hui Li²
Leslie Pack Kaelbling¹ Nicole Carey²

Abstract

Domain randomization in reinforcement learning is an established technique for increasing the robustness of control policies trained in simulation. By randomizing environment properties during training, the learned policy can become robust to uncertainties along the randomized dimensions. While the environment distribution is typically specified by hand, in this paper we investigate automatically discovering a sampling distribution via entropy-regularized reward maximization of a normalizing-flow-based neural sampling distribution. We show that this architecture is more flexible and provides greater robustness than existing approaches that learn simpler, parameterized sampling distributions, as demonstrated in six simulated and one real-world robotics domain. Lastly, we explore how these learned sampling distributions, along with a privileged value function, can be used for out-of-distribution detection in an uncertainty-aware multi-step manipulation planner.

1. Introduction

Reinforcement learning (RL) is a powerful tool in robotics because it can be used to learn effective control policies for systems with highly complex dynamics that are difficult to model analytically. Unlike traditional control methods, which rely on precise mathematical models, RL learns directly from simulated or real-world experience (Luo & Li, 2021; Zhu et al., 2020; Schoettler et al., 2020).

However, RL approaches can be inefficient, involving slow, minimally parallelized, and potentially unsafe data-gathering processes when performed in real environ-

ments (Kober et al., 2013). Learning in simulation eliminates some of these problems, but introduces new issues in the form of discrepancies between the training and real-world environments (Valassakis et al., 2020).

Successful RL from simulation hence requires efficient and accurate models of both robot and environment during the training process. But even with highly accurate geometric and dynamic simulators, the system can still be only considered partially observable (Kober et al., 2013)—material qualities, inertial properties, perception noise, contact and force sensor noise, manufacturing deviations and tolerances, and imprecision in robot calibration all add uncertainty to the model.

To improve the robustness of learned policies against sim-to-real discrepancies, it is common to employ domain randomization, varying the large set of environmental parameters inherent to a task according to a given underlying distribution (Muratore et al., 2019). In this way, policies are trained to maximize their overall performance over a diverse set of models. These sampling distributions are typically constructed manually with Gaussian or uniform distributions on individual parameters with hand-selected variances and bounds. However, choosing appropriate distributions for each of the domain randomization parameters remains a delicate process (Josifovski et al., 2022); too broad a distribution leads to suboptimal local minima convergence (see Figure 3), while too narrow a distribution leads to poor real-world generalization (Mozifian et al., 2019; Packer et al., 2018). Many existing methods rely on real-world rollouts from hardware experiments to estimate dynamics parameters (Chebotar et al., 2019; Ramos et al., 2019; Muratore et al., 2022). However, for complex tasks with physical parameters that are difficult to efficiently or effectively sample, this data may be time-consuming to produce, or simply unavailable.

An alternative strategy is to learn an environment distribution during training with the aim of finding the broadest possible training distribution that can feasibly be solved in order to maximize the chances of transferring to an unknown target environment. Automating updates to parameter distributions during the training process can remove the need

¹MIT CSAIL ²Autodesk Research. Correspondence to: Aidan Curtis <curtisa@mit.edu>.

for heuristic tuning and iterative experimentation (Mozifian et al., 2019; OpenAI et al., 2019; Tiboni et al., 2024). In this paper, we present GoFlow, a novel approach for learned domain randomization that combines actor-critic reinforcement learning architectures (Schulman et al., 2017; Haarnoja et al., 2018) with a neural sampling distribution to learn robust policies that generalize to real-world settings. By maximizing the diversity of parameters during sampling, we actively discover environments that are challenging for the current policy but still solvable given enough training.

As proof of concept, we investigate one real-world use case: contact-rich manipulation for assembly. Assembly is a critical area of research for robotics, requiring a diverse set of high-contact interactions which often involve wide force bandwidths and unpredictable dynamic changes. Recently, sim-to-real RL has emerged as a potentially useful strategy for learning robust contact-rich policies without laborious real-world interactions (Noseworthy et al., 2024; Tang et al., 2023b; Zhang et al., 2024). We build on this work by testing our method on the real-world industrial assembly task of gear insertion.

Lastly, we extend this classical gear insertion task to the setting of multi-step decision making under uncertainty and partial observability. As shown in this paper and elsewhere (Tiboni et al., 2024; Mozifian et al., 2019; OpenAI et al., 2019), policies trained in simulation have an upper bound on the environmental uncertainties that they can be conformant to. For example, a visionless robot executing an insertion policy can only tolerate so much in-hand pose error. However, estimates of this uncertainty can be used to inform high level control decisions through task-oriented information gathering (Liang et al., 2020; Curtis et al., 2023). For example, looking closer at objects can result in more accurate pose estimates or tracking objects in the hand to detect slippage. By integrating a probabilistic pose estimation model, we can use the sampling distributions learned with GoFlow as an out-of-distribution detector to determine whether the policy is expected to succeed under its current belief about the world state. If the robot has insufficient information, it can act to deliberately seek the needed information using a simple belief-space planning algorithm.

Our contributions are as follows: We introduce GoFlow, a novel domain randomization method that combines actor-critic reinforcement learning with a learned neural sampling distribution. We show that GoFlow achieves higher domain coverage than fixed and other learning-based solutions to domain randomization on a suite of simulated environments. We demonstrate the efficacy of GoFlow in a real-world contact-rich manipulation task—gear insertion—and extend it to multi-step decision-making under uncertainty. By integrating a probabilistic pose estimation model, we enable the robot to actively gather additional information when needed,

enhancing performance in partially observable settings.

2. Related Work

Recent developments in reinforcement learning have demonstrated that policies trained in simulation can effectively transfer to real-world robots, even for contact-rich robotic tasks (Zhang et al., 2024; Tang et al., 2023b; Noseworthy et al., 2024; Jin et al., 2023). A key innovation enabling this transferability is domain randomization (Peng et al., 2017), where environment parameters are sampled from predefined distributions during training, enabling learned policies to generalize to environmental uncertainties upon deployment.

Traditionally, domain randomization requires manually defining these sampling distributions, which can be labor-intensive. To address this, recent work has explored methods to automatically learn these distributions, aiming for maximal generalization with minimal manual effort. One notable approach involves constructing adversarial distributions that challenge the current policy, ensuring coverage of the environment parameter space (Mehta et al., 2020; Wang et al., 2025). While adversarial methods can outperform uniform randomization, they typically assume all environments sampled are solvable, limiting their effectiveness in highly uncertain scenarios.

An alternative line of research adopts a curriculum-based approach, progressively expanding the complexity of the training distribution while maintaining policy success. Specific curricular methods include minimizing divergence from a target distribution using multivariate Gaussians (Mozifian et al., 2019), maximizing entropy through independent beta distributions (Tiboni et al., 2024), and incrementally expanding uniform sampling distributions via boundary sampling (OpenAI et al., 2019). In this work, we propose a novel learned domain randomization technique leveraging normalizing flows (Rezende & Mohamed, 2015) as neural sampling distributions. This approach offers increased flexibility and expressivity over existing parametric methods.

Beyond training robust policies in simulation, learned sampling distributions can be tied to the real-world environmental conditions under which policies are likely to succeed. Previous works have integrated domain randomization with real-world interactions for more informed training distributions (Ramos et al., 2019; Sagawa & Hino, 2025; Ajay et al., 2022; Muratore et al., 2020) or to find the maximally effective real-world strategy (Yu et al., 2018; Ren et al., 2023). However, these methods often necessitate expensive policy retraining or data-intensive evolutionary search based on real-world feedback, posing challenges for real-time applications. Instead, we utilize our learned sampling distribution as an out-of-distribution detector within a multi-step planning framework, enabling fast and data-efficient

information gathering in the real world.

3. Background

3.1. Markov Decision Process

A Markov Decision Process (MDP) is a mathematical framework for modeling decision-making. Formally, an MDP is defined as a tuple $(\mathcal{S}, \mathcal{A}, P, R, \gamma)$, where \mathcal{S} is the state space, \mathcal{A} is the action space, $P : \mathcal{S} \times \mathcal{A} \times \mathcal{S} \rightarrow [0, 1]$ is the state transition probability function, where $P(s' | s, a)$ denotes the probability of transitioning to state s' from state s after taking action a , $R : \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}$ is the reward function, where $R(s, a)$ denotes the expected immediate reward received after taking action a in state s , $\gamma \in [0, 1]$ is the discount factor, representing the importance of future rewards.

A policy $\pi : \mathcal{S} \times \mathcal{A} \rightarrow [0, 1]$ defines a probability distribution over actions given states, where $\pi(a | s)$ is the probability of taking action a in state s . The goal is to find an optimal policy π^* that maximizes the expected cumulative discounted reward.

3.2. Domain Randomization

Domain randomization introduces variability into the environment by randomizing certain parameters during training. Let Ξ denote the space of domain randomization parameters, and let $\xi \in \Xi$ be a specific instance of these parameters. Each ξ corresponds to a different environment configuration or dynamics.

We can define a parameterized family of Markov Decision Processes (MDPs) where each $\mathcal{M}_\xi = (\mathcal{S}, \mathcal{A}, P_\xi, R_\xi, \gamma)$ has transition dynamics P_ξ and reward function R_ξ dependent on ξ . The agent interacts with environments sampled from a distribution over Ξ , typically denoted as $p(\xi)$.¹

The objective is to learn a policy $\pi : \mathcal{S} \rightarrow \mathcal{A}$ that maximizes the expected return across the distribution environments:

$$J(\pi) = \mathbb{E}_{\xi \sim p(\xi)} \left[\mathbb{E}_{\tau \sim P_{\xi, \pi}} \left[\sum_{t=0}^{\infty} \gamma^t R_\xi(s_t, a_t) \right] \right], \quad (1)$$

where $\tau = \{(s_0, a_0, s_1, a_1, \dots)\}$ denotes a trajectory generated by policy π in environment ξ . Domain randomization aims to find a policy π^* such that: $\pi^* = \arg \max_{\pi} J(\pi)$.

In deep reinforcement learning, the policy π is a neural network parameterized by θ , denoted as π_θ . The agent learns the policy parameters θ through interactions with simulated

¹This problem can also be thought of as a POMDP where the observation space is \mathcal{S} and the state space is a product of \mathcal{S} and Ξ as discussed in (Kwon et al., 2021).

environments sampled from $p(\xi)$. In our implementation, we employ the Proximal Policy Optimization (PPO) algorithm (Schulman et al., 2017), an on-policy policy gradient method that optimizes a stochastic policy while ensuring stable and efficient learning.

To further stabilize training, we pass privileged information about the environment parameters ξ to the critic network. The critic network, parameterized by ψ , estimates the state-value function:

$$V_\psi(s_t, \xi) = \mathbb{E}_{\pi_\theta} \left[\sum_{k=0}^{\infty} \gamma^k r_{t+k} \mid s_t, \xi \right], \quad (2)$$

where s_t is the current state, r_{t+k} are future rewards, and γ is the discount factor. By incorporating ξ , the critic can provide more accurate value estimates with lower variance (Pinto et al., 2018). The actor network $\pi_\theta(a_t | s_t)$ does not have access to ξ , ensuring that the policy relies only on observable aspects of the state.

3.3. Normalizing Flows

Normalizing flows are a class of generative models that transform a simple base distribution into a complex target distribution using a sequence of invertible, differentiable functions. Let $z \sim p_Z(z)$ be a latent variable from a base distribution (e.g., a standard normal distribution). A normalizing flow defines an invertible transformation $f_\phi : \mathbb{R}^d \rightarrow \mathbb{R}^d$ parameterized by neural network parameters ϕ , such that $x = f_\phi(z)$, aiming for x to follow the target distribution.

The density of x is computed using the change of variables formula:

$$p_X(x) = p_Z(f_\phi^{-1}(x)) \left| \det \left(\frac{\partial f_\phi^{-1}(x)}{\partial x} \right) \right|. \quad (3)$$

For practical computation, this is often rewritten as:

$$\log p_X(x) = \log p_Z(z) - \log \left| \det \left(\frac{\partial f_\phi(z)}{\partial z} \right) \right|, \quad (4)$$

where $\frac{\partial f_\phi(z)}{\partial z}$ is the Jacobian of f_ϕ at z . By composing multiple such transformations $f_\phi = f_{\phi_K} \circ \dots \circ f_{\phi_1}$, each parameterized by neural network parameters ϕ_k , normalizing flows can model highly complex distributions.

In our work, we employ *neural spline flows* (Durkan et al., 2019), a type of normalizing flow where the invertible transformations are constructed using spline-based functions. Specifically, the parameters ϕ represent the coefficients of

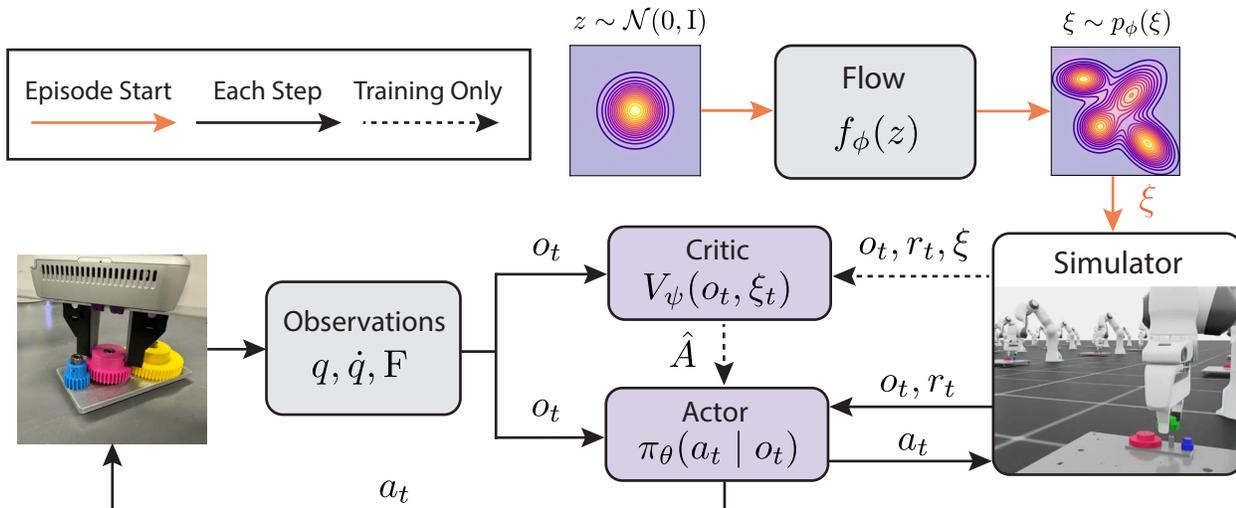


Figure 1. An architecture diagram for our actor-critic RL training setup in the gears domain. At the start of each training episode, an environment initialization $\xi \sim p_\phi(\xi)$ is sampled via the normalizing flow model. This variable conditions both the simulator and the privileged critic. At each time step, the agent receives observations including joint positions q , joint velocities \dot{q} , and force measurements F from which the actor policy generates actions. The critic takes observations o_t , reward r_t , and environment initialization ξ to estimate the value. This estimated value is used to compute advantage \hat{A} which is then used to train the actor. Only the actor is used during testing and real-world deployment, but the critic and flow can help determine when a policy will be successful as discussed in section 6.

the splines (e.g., knot positions and heights) and the weights and biases of the neural networks that parameterize these splines.

4. Method

In this section, we introduce **GoFlow**, a method for learned domain randomization that *goes with the flow* by adaptively adjusting the domain randomization process using normalizing flows.

In traditional domain randomization setups, the distribution $p(\xi)$ is predefined. However, selecting an appropriate $p(\xi)$ is crucial for the policy’s performance and generalization. Too broad a sampling distribution and the training focuses on unsolvable environments and falls into local minima. In contrast, too narrow a sampling distribution leads to poor generalization and robustness. Additionally, rapid changes to the sampling distribution can lead to unstable training. To address these challenges, prior works such as (Klink et al., 2021) have proposed a self-paced learner, which starts by mastering a small set of environments, and gradually expands the tasks to solver harder and harder problems while maintaining training stability. This strategy has subsequently been applied to domain randomization in (Mozifian et al., 2019) and (Tiboni et al., 2024), where terms were included for encouraging spread over the sampling

space for greater generalization. We take inspiration from these works to form a joint optimization problem:

$$\max_{p, \pi} \{ \mathbb{E}_{\xi \sim p} [J_\xi(\pi)] + \alpha \mathcal{H}(p) - \beta D_{KL}(p_{\text{old}} \| p) \} \quad (5)$$

where p is the sampling distribution over environment parameters $p(\xi)$, $\mathcal{H}(p)$ is the differential entropy of $p(\xi)$, $D_{KL}(p \| p_{\text{old}})$ is the divergence between the current and previous sampling distributions, and $\alpha > 0, \beta > 0$ are regularization coefficients that control the trade-off between generalizability, training stability, and the expected reward under the sampling distribution.

Other learned domain randomization approaches propose similar objectives. (Mozifian et al., 2019) maximizes reward but replaces entropy regularization with a KL divergence to a fixed target distribution and omits the self-paced KL term. (Tiboni et al., 2024) includes all three objectives but frames the reward and self-paced KL terms as constraints, maximizing entropy through a nonlinear optimization process that is not easily adaptable to neural sampling distributions. We compare GoFlow to these methods in our experiments to highlight its advantages. To our knowledge, GoFlow is the first method to optimize such an objective with a neural sampling distribution.

The GoFlow algorithm (Algorithm 1) begins by initial-

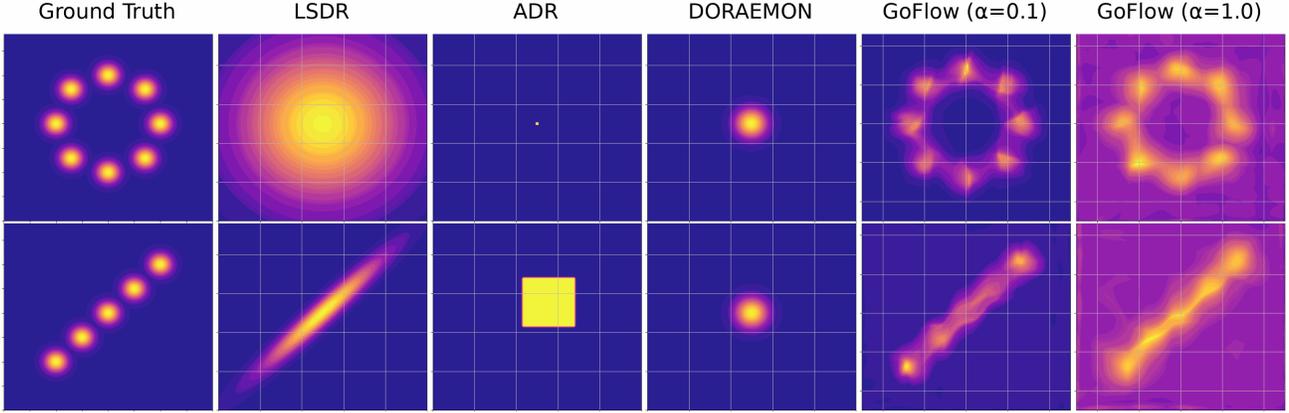


Figure 2. An illustrative domain showing the learned sampling functions over the space of unobserved parameters for the tested baselines. Compared to other learning methods, GoFlow correctly models the multimodality and inter-variable dependencies of the underlying reward function. This toy domain, along with other domains in our experiments, violates some of the assumptions made by prior works, such as the feasibility of the center point of the range.

Algorithm 1 GoFlow

Require: Initial policy parameters θ , flow parameters ϕ , training steps N , network updates K , entropy coefficient α , similarity coefficient β , and learning rate η_ϕ

- 1: **for** $n = 1$ to N **do**
- 2: Sample $\{\xi_i^{\text{train}}\}_{i=1}^B \sim p_\phi(\xi)$, $\{\xi_i^{\text{test}}\}_{i=1}^B \sim u(\xi)$
- 3: Train π_θ with ξ_i^{train} initializations
- 4: Estimate $J_{\xi_i^{\text{test}}}(\pi_\theta)$ via policy rollouts
- 5: Save current flow distribution as $p_{\phi_{\text{old}}}(\xi)$
- 6: **for** $k = 1$ to K **do**
- 7: $\mathcal{R} \leftarrow \frac{|\Xi|}{B} \sum_{i=1}^B \left[p_\phi(\xi_i^{\text{test}}) J_{\xi_i^{\text{test}}}(\pi_\theta) \right]$.
- 8: $\hat{H} \leftarrow -|\Xi| \cdot \mathbb{E}_{\xi \sim u(\xi)} \left[p_\phi(\xi) \log p_\phi(\xi) \right]$
- 9: $\hat{D}_{KL} \leftarrow \mathbb{E}_{\xi \sim p_{\phi_{\text{old}}}(\xi)} \left[\log p_{\phi_{\text{old}}}(\xi) - \log p_\phi(\xi) \right]$
- 10: $\phi \leftarrow \phi + \eta_\phi \nabla_\phi \left(\mathcal{R} + \alpha \hat{H} - \beta \hat{D}_{KL} \right)$
- 11: **end for**
- 12: **end for**

izing both the policy parameters θ and the normalizing flow parameters ϕ . In each training iteration, GoFlow first samples a batch of environment parameters $\{\xi_i\}_{i=1}^B$ from the current distribution modeled by the normalizing flow (Line 2). These sampled parameters are used to train the policy π_θ (Line 3). Following the policy update, expected returns $J_{\xi_i}(\pi_\theta)$ are estimated for each sampled environment through policy rollouts, providing a measure of the policy’s performance under a target uniform distribution $u(\xi)$ after being trained on $p(\xi)$ (Line 4).

After sampling these rollouts, GoFlow then performs K steps of optimization on the sampling distribution. GoFlow

first estimates the policy’s performance on the sampling distribution using the previously sampled rollout trajectories (Line 7). The entropy of the sampling distribution (Line 8) and divergence from the previous sampling distribution (Line 9) are estimated using newly drawn samples. Importantly, we compute the reward and entropy terms by importance sampling from a uniform distribution rather than from the flow itself. This broad coverage helps prevent the learned distribution from collapsing around a small region of parameter space. Derivations for these equations can be found in Appendix A.2. These terms are combined to form a loss, which is differentiated to update the parameters of the sampling distribution (Line 10). An architecture diagram for this approach can be seen in Figure 1.

5. Domain Randomization Experiments

Our simulated experiments compare policy robustness in a range of domains. For full details on the randomization parameters and bounds, see Appendix A.3.

5.1. Domains

First, we examine the application of GoFlow to an illustrative 2D domain that is multimodal and contains intervariable dependencies. The state and action space are in \mathbb{R}^2 . The agent is initialized randomly in a bounded x, y plane. An energy function is defined by a composition of Gaussians placed in a regular circular or linear array. The agent can observe its position with Gaussian noise proportional to the inverse of the energy function. The agent is rewarded for guessing its location, but is incapable of moving. This task is infeasible when the agent is sufficiently far from any

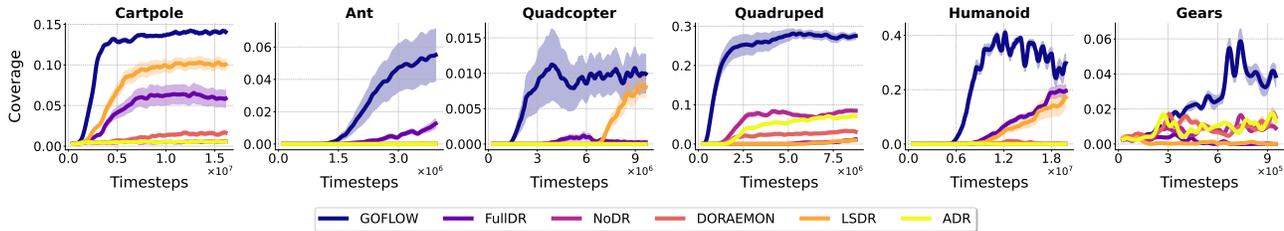


Figure 3. The coverage ratio over the target distribution across five random seeds for each of the environments. The bands around each curve indicate the standard error.

of the Gaussian centers, so a sampling distribution should come to resemble the energy function. Some example functions learned by GoFlow and baselines from Section 5.2 can be seen in Figure 2.

Second, we quantitatively compare GoFlow to existing baselines including Cartpole, Ant, Quadcopter, Quadruped, and Humanoid in the IsaacLab suite of environments (Mittal et al., 2023). We randomize over parameters such as link masses, joint frictions, and material properties.

Lastly, we evaluate our method on a contact-rich robot manipulation task of gear insertion, a particularly relevant problem for robotic assembly. In the gears domain, we randomize over the relative pose between the gripper and the held gears along three degrees of freedom. The problem is made difficult by the uncertainty the robot has about the precise location of the gear relative to the hand. The agent must learn to rely on signals of proprioception and force feedback to guide the gear into the gear shaft. The action space consists of end-effector pose offsets along three translational degrees of freedom and one rotational degree around the z dimension. The observation space consists of a history of the ten previous end effector poses and velocities estimated via finite differencing. In addition to simulated experiments, our trained policies are tested on a Franka Emika robot using IndustReal library built on the frankapy toolkit (Tang et al., 2023a; Zhang et al., 2020).

5.2. Baselines

In our domain randomization experiments, we compare to a number of standard RL baselines and learning-based approaches from the literature. In our quantitative experiments, success is determined by the sampled environment passing a certain performance threshold J_T that was selected for each environment. We measure task performance in terms of *coverage*, which is defined as the proportion of the total sampling distribution for which the policy receives higher than J_t reward. Coverage is estimated via policy rollouts on 4096 uniformly selected environment initialization samples from the total range of parameters. All baselines are trained with

an identical neural architecture and PPO implementation. The success thresholds along with other hyperparameters are in Appendix A.5.

We evaluate on the following baselines. First, we compare to no domain randomization (**NoDR**) which trains on a fixed environment parameter at the centroid of the parameter space. Next, we compare to a full domain randomization (**FullDR**) which samples uniformly across the domain within the boundaries during training. In addition to these fixed randomization methods, we evaluate against some other learning-based solutions from the literature: **ADR** (OpenAI et al., 2019) learns uniform intervals that expand over time via boundary sampling. It starts by occupying an initial percentage of the domain and performs “boundary sampling” during training with some probability. The rewards attained from boundary sampling are compared to thresholds that determine if the boundary should be expanded or contracted. **LSDR** (Mozifian et al., 2019) learns a multivariate gaussian sampling distribution using reward maximization with a KL divergence regularization term weighted by an α hyperparameter. Lastly, **DORAE-MON** (Tiboni et al., 2024) learns independent beta distributions for each dimension of the domain, using a maximum entropy objective constrained by an estimated success rate.

We compare coverage across environments sampled from a uniform testing distribution within the environment bounds. Our findings in Figure 3 show that GoFlow matches or outperforms baselines across all domains. We find that our method performs particularly well in comparison to other learned baselines when the simpler or more feasible regions of the domain are off-center, irregularly shaped, and have inter-parameter dependencies such as those seen in Figure 2. We intentionally chose large intervals with low coverage to demonstrate this capability, and we perform a full study of how baselines degrade significantly with increased parameter ranges while GoFlow degrades more gracefully (see Appendix A.6). Additionally, coverage should not be mistaken for success-rate, as these policies can be integrated into a planning framework that avoids the use of infeasible actions or gathers information to increase coverage as

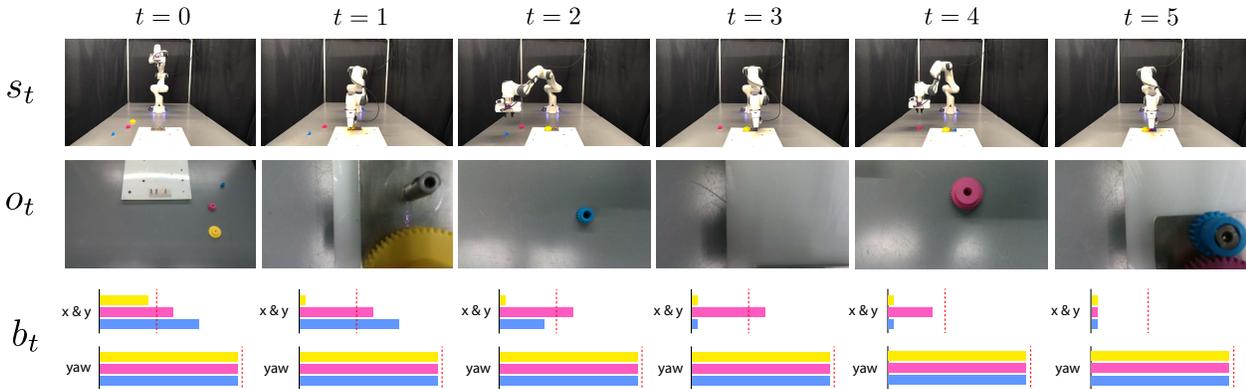


Figure 4. A multi-step manipulation plan using probabilistic pose estimation to estimate and update beliefs over time. The three rows show the robot state s_t , the observation o_t , and the robot belief b_t at each timestep. The red dotted line in the belief indicates the marginal entropy thresholds for the x, y, and yaw (rotation around z) dimensions as determined by the learned normalizing flow. A belief with entropy surpassing the threshold line indicates the policy will likely fail. For full visualizations of the belief posteriors, flow distributions, and value maps, see Figure 14.

discussed in the following sections.

In addition to simulated experiments, we additionally evaluated the trained policies on a real-world gear insertion task. The results of those real-world experiments show that GoFlow results in more robust sim-to-real transfer as seen in Table 1 and in the supplementary videos.

6. Application to Multi-step manipulation

While reinforcement learning has proven to be a valuable technique for learning short-horizon skills in dynamic and contact-rich settings, it often struggles to generalize to more long-horizon and open ended problems (Sutton & Barto, 2018). The topic of sequencing short horizon skills in the context of a higher-level decision strategy has been of increasing interest to both the planning (Mishra et al., 2023) and reinforcement learning communities (Nasiriany et al., 2021). For this reason, we examine the utility of these learned sampling distributions as out-of-distribution detectors, or belief-space preconditions, in the context of a multi-step planning system.

6.1. Belief-space planning background

Belief-space planning is a framework for decision-making under uncertainty, where the agent maintains a probability distribution over possible states, known as the *belief state*. Instead of planning solely in the state space \mathcal{S} , the agent operates in the *belief space* \mathcal{B} , which consists of all possible probability distributions over \mathcal{S} . This approach is particularly useful in partially observable environments where there is uncertainty in environment parameters and where it

is important to take actions to gain information.

Rather than operating at the primitive action level, belief-space planners often make use of high-level actions \mathcal{A}_Π , sometimes called skills or options. In our case, these high-level actions will be a discrete set of pretrained RL policies. These high-level actions come with a *belief-space precondition* and a *belief-space effect*, both of which are subsets of the belief space \mathcal{B} (Kaelbling & Lozano-Perez, 2013; Curtis et al., 2024). Specifically, a high-level action $\pi \in \mathcal{A}_\Pi$ is associated with two components: a precondition $\text{Pre}_\pi \subseteq \mathcal{B}$, representing the set of belief states from which the action can be applied, and an effect $\text{Eff}_\pi \subseteq \mathcal{B}$, representing the set of belief states that the action was designed or trained to achieve. If the precondition holds—that is, if the current belief state b satisfies $b \in \text{Pre}_\pi$ —then applying action a will achieve the effect Eff_π with probability at least $\eta \in [0, 1]$. More formally, if $b \in \text{Pre}_\pi$, then $\Pr(b_{t+1} \in \text{Eff}_\pi \mid b_t, \pi) \geq \eta$. Depending on the planner, η may be set in advance, or calculated by the planner as a function of the belief.

This formalization allows planners to reason abstractly about the effects of high-level actions under uncertainty, which can result in generalizable decision-making in long-horizon problems that require active information gathering or risk-awareness.

6.2. Computing preconditions

In this section, we highlight the potential application of the learned sampling distribution p_ϕ and privileged value function V_ψ as a useful artifacts for belief-space planning. In particular, we are interested in identifying belief-space

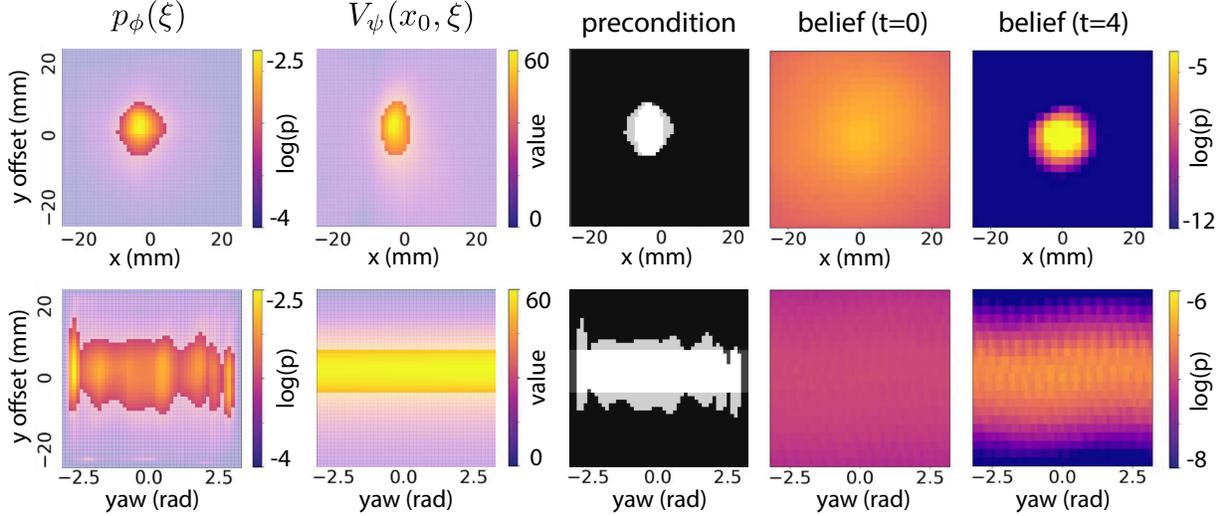


Figure 5. A visual example of the precondition computation described in Section 6.2 for the gear assembly plan shown in Figure 4. The two rows show two different projections of the 3D sampling space (x position vs y position in the top row and y position vs yaw rotation in the bottom row). We apply a threshold ϵ to the sampling distribution to remove low-probability regions (column 1). Additionally, we filter the value function by retaining only the regions where the expected value exceeds a predetermined threshold η (column 2). The intersection of these two regions defines the belief-space precondition, indicating where the policy is likely to succeed (column 3). Comparing the precondition to the beliefs, we can see that the belief is not sufficiently contained within the precondition at $t = 0$ (column 4), but passes the success threshold η at after closer inspection at $t = 4$ (column 5).

preconditions of a set of trained skills.

One point of leverage we have for this problem is the privileged value function $V_\psi(s, \xi)$, which was learned alongside the policy during training. One way to estimate the belief-space precondition is to simply find the set of belief states for which the expected value of the policy is larger than J_T with probability greater than η under the belief:

$$\text{Pre}_\pi = \{b \in \mathcal{B} \mid \mathbb{E}_b [\mathbb{1}_{V_\psi(s, \xi) > J_T}] > \eta\}. \quad (6)$$

However, a practical issue with this computation is that the value function is likely not calibrated in large portions of the state space that were not seen during policy training. To address this, we focus on regions of the environment where the agent has higher confidence due to sufficient sampling, i.e., where $p_\phi(\xi) > \epsilon$ for a threshold ϵ . This enables us to integrate the value function over the belief distribution $b(x, \xi)$ and the trusted region within Ξ :

$$\text{Pre}_\pi = \{b \in \mathcal{B} \mid \mathbb{E}_b [\mathbb{1}_{V_\psi(s, \xi) > J_T} \cdot \mathbb{1}_{p_\phi(\xi) > \epsilon}] > \eta\} \quad (7)$$

Here, η lower bounds the probability of achieving the desired effects Eff_π (or value greater than J_T) after executing π in any belief state in Pre_π . Figure 5 shows an example precondition for a single step of the assembly plan.

6.3. Updating beliefs

Updating the belief state requires a probabilistic state estimation system that outputs a posterior over the unobserved environment variables, rather than a single point estimate. We use a probabilistic object pose estimation framework called Bayes3D to infer posterior distributions over object pose (Gothoskar et al., 2023). For details on this, see Appendix A.4.1.

The benefit of this approach in contrast to traditional rendering-based pose estimation systems, such as those presented in (Wen et al., 2024) or (Labbé et al., 2022), is that pose estimates from Bayes3D indicate high uncertainty for distant, small, or occluded objects as well as uncertainty stemming from object symmetry. Figure 14 shows the pose beliefs across the multi-step plan.

6.4. A simple belief-space planner

While the problem of general-purpose multi-step planning in belief-space has been widely studied, in this paper we use a simple BFS belief-space planner to demonstrate the utility of the learned sampling distributions as belief-space preconditions. The full algorithm can be found in Algorithm 2.

An example plan can be seen in Figure 4. The goal is to assemble the gear box by inserting all three gears (yellow, pink, and blue) into the shafts on the gear plate. Each gear insertion is associated with a separate policy for each color

trained with GoFlow. In addition to the trained policies, the robot is given access to an object-parameterized inspection action which has no preconditions and whose effects are a reduced-variance pose estimate attained by moving the camera closer to the object. The robot is initially uncertain of the x , y , and yaw components of the 6-dof pose based on probabilistic pose estimates. Despite this uncertainty, the robot is confident enough in the pose of the largest and closest yellow gear to pick it up and insert it. In contrast, the blue and pink gears require further inspection to get a better pose estimate. Closer inspection reduces uncertainty along the x and y axis, but reveals no additional information about yaw dimension due to rotational symmetry. Despite an unknown yaw dimension, the robot is confident in the insertion because the flow p_ϕ indicates that success is invariant to the yaw dimension. This is due to the fact that success in the insertion task is defined by the distance between the bottom center of the gear and the base of the gear shaft, which is independent of gear rotation. For visualizations of the beliefs and flows at each step, see Appendix A.4.

7. Conclusion and discussion

In this paper, we introduced GoFlow, a novel approach to domain randomization that uses normalizing flows to dynamically adjust the sampling distribution during reinforcement learning. By combining actor-critic reinforcement learning with a learned neural sampling distribution, we enabled more flexible and expressive parameterization of environmental variables, leading to better generalization in complex tasks like contact-rich assembly. Our experiments demonstrated that GoFlow achieves higher coverage than traditional fixed and learning-based domain randomization techniques across a variety of simulated environments, particularly in scenarios where the domain has irregular dependencies between parameters. The method also showed promise in real-world robotic tasks including contact-rich assembly.

Moreover, we extended GoFlow to multi-step decision-making tasks, integrating it with belief-space planning to handle long-horizon problems under uncertainty. This extension enabled the use of learned sampling distributions and value functions as preconditions leading to active information gathering.

Although GoFlow enables more expressive sampling distributions, it also presents some new challenges. One limitation of our method is that it has higher variance due to occasional training instability of the flow. This instability can be alleviated by increasing β , but at the cost of reduced sample efficiency (see Appendix A.5). In addition, using the flow and value estimates for belief-space planning require manual selection and tuning of several thresholds which are environment specific. The η parameter may be converted

from a threshold into a cost in the belief space planner, which would remove one point of manual tuning. However, removing the ϵ parameter may prove more difficult, as it would require uncertainty quantification of the neural value function. Despite these challenges, we hope this work inspires further research on integrating short-horizon learned policies into broader planning frameworks, particularly in contexts involving uncertainty and partial observability.

Acknowledgements

We gratefully acknowledge support from NSF grant 2214177; from AFOSR grant FA9550-22-1-0249; from ONR MURI grants N00014-22-1-2740 and N00014-24-1-2603; from the MIT Quest for Intelligence; and from the Robotics and Artificial Intelligence Institute. Aidan Curtis was supported by the NSF Graduate Research Fellowship. This work was partially conducted during a summer internship at Autodesk.

Impact Statement

This paper presents work whose goal is to advance the field of Machine Learning. There are many potential societal consequences of our work, none which we feel must be specifically highlighted here.

References

- Ajay, A., Gupta, A., Ghosh, D., Levine, S., and Agrawal, P. Distributionally adaptive meta reinforcement learning. In *Proceedings of the 36th International Conference on Neural Information Processing Systems, NIPS '22*, Red Hook, NY, USA, 2022. Curran Associates Inc. ISBN 9781713871088.
- Chebatar, Y., Handa, A., Makoviychuk, V., Macklin, M., Issac, J., Ratliff, N., and Fox, D. Closing the sim-to-real loop: Adapting simulation randomization with real world experience. In *2019 International Conference on Robotics and Automation (ICRA)*, pp. 8973–8979. IEEE, 2019.
- Curtis, A., Kaelbling, L., and Jain, S. Task-directed exploration in continuous pomdps for robotic manipulation of articulated objects. In *2023 IEEE International Conference on Robotics and Automation (ICRA)*, pp. 3721–3728, 2023. doi: 10.1109/ICRA48891.2023.10160306.
- Curtis, A., Matheos, G., Gothoskar, N., Mansinghka, V., Tenenbaum, J., Lozano-Pérez, T., and Kaelbling, L. P. Partially observable task and motion planning with uncertainty and risk awareness, 2024. URL <https://arxiv.org/abs/2403.10454>.
- Del Moral, P., Doucet, A., and Jasra, A. Sequential monte

- carlo samplers. *Journal of the Royal Statistical Society Series B: Statistical Methodology*, 68(3):411–436, 2006.
- Durkan, C., Bekasov, A., Murray, I., and Papamakarios, G. Neural spline flows. In *Advances in Neural Information Processing Systems (NeurIPS)*, volume 32, 2019.
- Gothoskar, N., Ghavami, M., Li, E., Curtis, A., Noseworthy, M., Chung, K., Patton, B., Freeman, W. T., Tenenbaum, J. B., Klukas, M., et al. Bayes3d: fast learning and inference in structured generative models of 3d objects and scenes. *arXiv preprint arXiv:2312.08715*, 2023.
- Haarnoja, T., Zhou, A., Abbeel, P., and Levine, S. Soft actor-critic: Off-policy maximum entropy deep reinforcement learning with a stochastic actor. In *International conference on machine learning*, pp. 1861–1870. PMLR, 2018.
- Jin, P., Lin, Y., Song, Y., Li, T., and Yang, W. Vision-force-fused curriculum learning for robotic contact-rich assembly tasks. *Frontiers in Neurorobotics*, 17:1280773, October 2023. doi: 10.3389/fnbot.2023.1280773.
- Josifovski, J., Malmir, M., Klarmann, N., Žagar, B. L., Navarro-Guerrero, N., and Knoll, A. Analysis of randomization effects on sim2real transfer in reinforcement learning for robotic manipulation tasks. In *2022 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, pp. 10193–10200. IEEE, 2022.
- Kaelbling, L. and Lozano-Perez, T. Integrated task and motion planning in belief space. *The International Journal of Robotics Research*, 32:1194–1227, 08 2013. doi: 10.1177/0278364913484072.
- Klink, P., Abdulsamad, H., Belousov, B., D’Eramo, C., Peters, J., and Pajarinen, J. A probabilistic interpretation of self-paced learning with applications to reinforcement learning. *J. Mach. Learn. Res.*, 22(1), January 2021. ISSN 1532-4435.
- Kober, J., Bagnell, J. A., and Peters, J. Reinforcement learning in robotics: A survey. *The International Journal of Robotics Research*, 32(11):1238–1274, 2013.
- Kwon, J., Efroni, Y., Caramanis, C., and Mannor, S. RL for latent mdps: regret guarantees and a lower bound. In *Proceedings of the 35th International Conference on Neural Information Processing Systems, NIPS ’21*, Red Hook, NY, USA, 2021. Curran Associates Inc. ISBN 9781713845393.
- Labbé, Y., Manuelli, L., Mousavian, A., Tyree, S., Birchfield, S., Tremblay, J., Carpentier, J., Aubry, M., Fox, D., and Sivic, J. Megapose: 6d pose estimation of novel objects via render compare, 2022. URL <https://arxiv.org/abs/2212.06870>.
- Liang, J., Saxena, S., and Kroemer, O. Learning active task-oriented exploration policies for bridging the sim-to-real gap. *CoRR*, abs/2006.01952, 2020. URL <https://arxiv.org/abs/2006.01952>.
- Luo, J. and Li, H. A learning approach to robot-agnostic force-guided high precision assembly. In *2021 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, pp. 2151–2157. IEEE, 2021.
- Mehta, B., Diaz, M., Golemo, F., Pal, C. J., and Paull, L. Active domain randomization. In *Conference on Robot Learning*, pp. 1162–1176. PMLR, 2020.
- Mishra, U. A., Xue, S., Chen, Y., and Xu, D. Generative skill chaining: Long-horizon skill planning with diffusion models, 2023. URL <https://arxiv.org/abs/2401.03360>.
- Mittal, M., Yu, C., Yu, Q., Liu, J., Rudin, N., Hoeller, D., Yuan, J. L., Singh, R., Guo, Y., Mazhar, H., Mandekar, A., Babich, B., State, G., Hutter, M., and Garg, A. Orbit: A unified simulation framework for interactive robot learning environments. *IEEE Robotics and Automation Letters*, 8(6):3740–3747, 2023. doi: 10.1109/LRA.2023.3270034.
- Mozifian, M., Higuera, J. C. G., Meger, D., and Dudek, G. Learning domain randomization distributions for transfer of locomotion policies. *CoRR*, abs/1906.00410, 2019. URL <http://arxiv.org/abs/1906.00410>.
- Muratore, F., Gienger, M., and Peters, J. Assessing transferability from simulation to reality for reinforcement learning. *IEEE transactions on pattern analysis and machine intelligence*, 43(4):1172–1183, 2019.
- Muratore, F., Eilers, C., Gienger, M., and Peters, J. Bayesian domain randomization for sim-to-real transfer. *CoRR*, abs/2003.02471, 2020. URL <https://arxiv.org/abs/2003.02471>.
- Muratore, F., Gruner, T., Wiese, F., Belousov, B., Gienger, M., and Peters, J. Neural posterior domain randomization. In Faust, A., Hsu, D., and Neumann, G. (eds.), *Proceedings of the 5th Conference on Robot Learning*, volume 164 of *Proceedings of Machine Learning Research*, pp. 1532–1542. PMLR, 08–11 Nov 2022.
- Nasiriany, S., Liu, H., and Zhu, Y. Augmenting reinforcement learning with behavior primitives for diverse manipulation tasks. *2022 International Conference on Robotics and Automation (ICRA)*, pp. 7477–7484, 2021.
- Noseworthy, M., Tang, B., Wen, B., Handa, A., Roy, N., Fox, D., Ramos, F., Narang, Y., and Akinola, I. Forge: Force-guided exploration for robust contact-rich manipulation under uncertainty, 2024. URL <https://arxiv.org/abs/2408.04587>.

- OpenAI, Akkaya, I., Andrychowicz, M., Chociej, M., Litwin, M., McGrew, B., Petron, A., Paino, A., Plappert, M., Powell, G., Ribas, R., Schneider, J., Tezak, N., Tworek, J., Welinder, P., Weng, L., Yuan, Q., Zaremba, W., and Zhang, L. Solving rubik’s cube with a robot hand. *CoRR*, abs/1910.07113, 2019. URL <http://arxiv.org/abs/1910.07113>.
- Packer, C., Gao, K., Kos, J., Krähenbühl, P., Koltun, V., and Song, D. Assessing generalization in deep reinforcement learning. *CoRR*, abs/1810.12282, 2018. URL <http://arxiv.org/abs/1810.12282>.
- Peng, X. B., Andrychowicz, M., Zaremba, W., and Abbeel, P. Sim-to-real transfer of robotic control with dynamics randomization. *CoRR*, abs/1710.06537, 2017. URL <http://arxiv.org/abs/1710.06537>.
- Pinto, L., Andrychowicz, M., Welinder, P., Zaremba, W., and Abbeel, P. Asymmetric actor critic for image-based robot learning. *RSS*, 2018.
- Ramos, F., Possas, R. C., and Fox, D. Bayessim: adaptive domain randomization via probabilistic inference for robotics simulators. *CoRR*, abs/1906.01728, 2019. URL <http://arxiv.org/abs/1906.01728>.
- Ren, A. Z., Dai, H., Burchfiel, B., and Majumdar, A. Adapt-sim: Task-driven simulation adaptation for sim-to-real transfer, 2023. URL <https://arxiv.org/abs/2302.04903>.
- Rezende, D. J. and Mohamed, S. Variational inference with normalizing flows. In *Proceedings of the 32nd International Conference on Machine Learning (ICML)*. PMLR, 2015.
- Rozet, F. et al. Zuko: Normalizing flows in pytorch, 2022. URL <https://pypi.org/project/zuko>.
- Sagawa, S. and Hino, H. Gradual domain adaptation via normalizing flows. *Neural Computation*, 37(3):522–568, 2025. doi: 10.1162/neco_a.01734.
- Schoettler, G., Nair, A., Luo, J., Bahl, S., Ojea, J. A., Solowjow, E., and Levine, S. Deep reinforcement learning for industrial insertion tasks with visual inputs and natural rewards. In *2020 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, pp. 5548–5555. IEEE, 2020.
- Schulman, J., Wolski, F., Dhariwal, P., Radford, A., and Klimov, O. Proximal policy optimization algorithms. *CoRR*, abs/1707.06347, 2017. URL <http://arxiv.org/abs/1707.06347>.
- Sutton, R. S. and Barto, A. G. *Reinforcement Learning: An Introduction*. A Bradford Book, Cambridge, MA, USA, 2018. ISBN 0262039249.
- Tang, B., Lin, M., Akinola, I., Handa, A., Sukhatme, G., Ramos, F., Fox, D., and Narang, Y. Industreal: Transferring contact-rich assembly tasks from simulation to reality. 07 2023a. doi: 10.15607/RSS.2023.XIX.039.
- Tang, B., Lin, M. A., Akinola, I., Handa, A., Sukhatme, G. S., Ramos, F., Fox, D., and Narang, Y. Industreal: Transferring contact-rich assembly tasks from simulation to reality, 2023b. URL <https://arxiv.org/abs/2305.17110>.
- Tiboni, G., Klink, P., Peters, J., Tommasi, T., D’Eramo, C., and Chalvatzaki, G. Domain randomization via entropy maximization, 2024. URL <https://arxiv.org/abs/2311.01885>.
- Valassakis, E., Ding, Z., and Johns, E. Crossing the gap: A deep dive into zero-shot sim-to-real transfer for dynamics. In *2020 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, pp. 5372–5379. IEEE, 2020.
- Wang, C., Lv, Y., Mao, Y., Qu, Y., Xu, Y., and Ji, X. Robust fast adaptation from adversarially explicit task distribution generation, 2025. URL <https://arxiv.org/abs/2407.19523>.
- Wen, B., Yang, W., Kautz, J., and Birchfield, S. Foundationpose: Unified 6d pose estimation and tracking of novel objects, 2024. URL <https://arxiv.org/abs/2312.08344>.
- Yu, W., Liu, C. K., and Turk, G. Policy transfer with strategy optimization. *CoRR*, abs/1810.05751, 2018. URL <http://arxiv.org/abs/1810.05751>.
- Zhang, K., Sharma, M., Liang, J., and Kroemer, O. A modular robotic arm control stack for research: Franka-interface and frankapy. *arXiv:2011.02398*, 2020.
- Zhang, X., Tomizuka, M., and Li, H. Bridging the sim-to-real gap with dynamic compliance tuning for industrial insertion. In *2024 IEEE International Conference on Robotics and Automation (ICRA)*, pp. 4356–4363, 2024. doi: 10.1109/ICRA57147.2024.10610707.
- Zhu, H., Yu, J., Gupta, A., Shah, D., Hartikainen, K., Singh, A., Kumar, V., and Levine, S. The ingredients of real-world robotic reinforcement learning. *arXiv preprint arXiv:2004.12570*, 2020.

A. Appendix

A.1. Code Release

The codebase for the project can be found [here](#).

A.2. Importance Sampling Proofs

Here we show how we obtain the importance-sampled estimates for both the reward term \mathcal{R} and the entropy term \mathcal{H} in GoFlow (Algorithm 1). Specifically, we sample from a uniform distribution $u(\xi)$ over Ξ , rather than from the learned distribution $p_\phi(\xi)$ directly, in order to avoid collapse onto narrow regions of the parameter space.

A.2.1. REWARD TERM: \mathcal{R}

We want an unbiased estimate of

$$\mathbb{E}_{\xi \sim p_\phi(\xi)} [J_\xi(\pi)] = \int_{\Xi} p_\phi(\xi) J_\xi(\pi) d\xi.$$

Let $u(\xi)$ be a uniform distribution over Ξ . Since $p_\phi(\xi) J_\xi(\pi) = \frac{p_\phi(\xi)}{u(\xi)} u(\xi) J_\xi(\pi)$, we can rewrite:

$$\int_{\Xi} p_\phi(\xi) J_\xi(\pi) d\xi = \int_{\Xi} u(\xi) \frac{p_\phi(\xi)}{u(\xi)} J_\xi(\pi) d\xi.$$

Hence, sampling $\xi_i \sim u(\xi)$ and averaging $\frac{p_\phi(\xi_i)}{u(\xi_i)} J_{\xi_i}(\pi)$ gives an unbiased Monte Carlo estimate of $\mathbb{E}_{p_\phi} [J_\xi(\pi)]$.

If Ξ has finite measure $|\Xi|$, then $u(\xi) = 1/|\Xi|$. Thus

$$\frac{p_\phi(\xi)}{u(\xi)} = p_\phi(\xi) |\Xi|.$$

Therefore, the empirical estimate becomes

$$\begin{aligned} \mathcal{R} &= \frac{1}{B} \sum_{i=1}^B [(p_\phi(\xi_i) |\Xi|) J_{\xi_i}(\pi)] \\ &= \frac{|\Xi|}{B} \sum_{i=1}^B p_\phi(\xi_i) J_{\xi_i}(\pi). \end{aligned}$$

where $\{\xi_i\}_{i=1}^B \sim u(\xi)$. This matches Line 7 in Algorithm 1.

A.2.2. ENTROPY TERM: \mathcal{H}

Similarly, to compute the differential entropy of $p_\phi(\xi)$, we have

$$\mathcal{H}(p_\phi) = - \int_{\Xi} p_\phi(\xi) \log p_\phi(\xi) d\xi.$$

Again, we apply the same importance sampling trick via $u(\xi)$. We write:

$$p_\phi(\xi) \log p_\phi(\xi) = \frac{p_\phi(\xi)}{u(\xi)} u(\xi) \log p_\phi(\xi).$$

Hence

$$\int_{\Xi} p_\phi(\xi) \log p_\phi(\xi) d\xi = \int_{\Xi} u(\xi) \frac{p_\phi(\xi)}{u(\xi)} \log p_\phi(\xi) d\xi.$$

If $|\Xi|$ is the measure of Ξ under $u(\xi)$, then $u(\xi) = 1/|\Xi|$. So the Monte Carlo estimate for $\int p_\phi(\xi) \log p_\phi(\xi) d\xi$ becomes

$$\frac{|\Xi|}{B} \sum_{i=1}^B p_\phi(\xi_i) \log p_\phi(\xi_i),$$

with $\xi_i \sim u(\xi)$. Multiplying by -1 yields the differential entropy:

$$\begin{aligned} \mathcal{H}(p_\phi) &= - \int_{\Xi} p_\phi(\xi) \log p_\phi(\xi) d\xi \\ &\approx - \frac{|\Xi|}{B} \sum_{i=1}^B p_\phi(\xi_i) \log p_\phi(\xi_i), \end{aligned}$$

which is exactly what we implement in Line 8 of Algorithm 1.

Remark A.1. Using uniform sampling $u(\xi)$ to approximate these terms provides global coverage of Ξ , helping prevent the learned distribution p_ϕ from collapsing around a small subset of parameter space. By contrast, if one sampled ξ from $p_\phi(\xi)$ itself for these terms, the distribution might fail to expand to other promising regions once it becomes peaked.

A.3. Domain Randomization Parameters

Below we describe the randomization ranges and parameter names for each environment. We also provide the reward success threshold (J_T) and cut the max duration of some environments in order to speed up training (t_{max}). J_t was chosen to be below the optimal performance under no environment randomization. We verified that the trained policy still exhibited qualitatively successful performance at the target reward threshold. Lastly, we slightly modified the Quadruped environment to only take a fixed forward command rather than the goal-conditioned policy learned by default. Other than those changes, the first four simulated environments official IsaacLab implementation.

- **Cartpole parameters** ($J_T = 50, t_{max} = 2s$):
 - Pole mass: Min = 0.01, Max = 20.0
 - Cart mass: Min = 0.01, Max = 20.0
 - Slider-Cart Friction: Min Bound = 0.0, Max Bound = 1.0
- **Ant parameters** ($J_T = 700, t_{max} = 2s$):
 - Torso mass: Min = 0.01, Max = 20.0

- **Quadcopter parameters** ($J_T = 15, t_{max} = 2s$):
 - Quadcopter mass: Min = 0.01, Max = 20.0
- **Quadruped parameters** ($J_T = 1.5, t_{max} = 5s$):
 - Body mass: Min = 0.0, Max = 200.0
 - Left front hip friction: Min = 0.0, Max = 0.1
 - Left back hip friction: Min = 0.0, Max = 0.1
 - Right front hip friction: Min = 0.0, Max = 0.1
 - Right back hip friction: Min = 0.0, Max = 0.1
- **Humanoid parameters** ($J_T = 1000, t_{max} = 5s$):
 - Torso Mass: Min = 0.01, Max = 25.0
 - Head Mass: Min = 0.01, Max = 25.0
 - Left Hand Mass: Min = 0.01, Max = 30.0
 - Right Hand Mass: Min = 0.01, Max = 30.0
- **Gear parameters** ($J_T = 50, t_{max} = 4s$):
 - Grasp Pose x: Min = -0.05, Max = 0.05
 - Grasp Pose y: Min = -0.05, Max = 0.05
 - Grasp Pose yaw: Min = -0.393, Max = 0.393

A.4. Multi-Step Planning Details

A.4.1. UPDATING BELIEFS VIA PROBABILISTIC POSE ESTIMATION

Updating the belief state b requires a probabilistic state estimation system that outputs a posterior over the state space S , rather than a single point estimate. Given a new observation o , we use a probabilistic object pose estimation framework (Bayes3D) to infer posterior distributions over object pose (Gothoskar et al., 2023).

The pose estimation system uses inference in an probabilistic generative graphics model with uniform priors on the translational x , y , and rotational yaw (or r_x) components of the 6-dof pose (since the object is assumed to be in flush contact with the table surface) and an image likelihood $P(o_{\text{rgbd}} | r_x, x, y)$. The object’s geometry and color information is given by a mesh model. The image likelihood is computed by rendering a latent image im^{rgbd} with the object pose corresponding to (r_x, x, y) and calculating the per-pixel likelihood:

$$P(o_{\text{rgbd}} | r_x, x, y) \propto \prod_{i,j \in C} \left[p_{\text{out}} + (1 - p_{\text{out}}) \cdot P_{\text{in}}(o_{i,j}^{\text{rgbd}} | r_x, x, y) \right] \quad (8)$$

$$P_{\text{in}}(o_{i,j}^{\text{rgbd}} | r_x, x, y) \propto \exp \left(- \frac{\|o_{i,j}^{\text{rgb}} - im_{i,j}^{\text{rgb}}\|_1}{b^{\text{rgb}}} - \frac{\|o_{i,j}^{\text{d}} - im_{i,j}^{\text{d}}\|_1}{b^{\text{d}}} \right) \quad (9)$$

where i and j are pixel row and column indices, C is the set of valid pixels returned by the renderer, b_{rgb} and b_{d} are hyperparameters that control posterior sensitivity to the color and depth channels, and p_{out} is the pixel outlier probability hyperparameter. For an observation o_{rgbd} , we can sample from $P(r_x, x, y | o_{\text{rgbd}}) \propto P(o_{\text{rgbd}} | r_x, x, y)$ to recover the object pose posterior with a tempering exponential factor α to encourage smoothness. We first find the maximum a posteriori (MAP) estimate of object pose using coarse-to-fine sequential Monte Carlo sampling (Del Moral et al., 2006) and then calculate a posterior approximation using a grid centered at the MAP estimate.

The benefit of this approach in contrast to traditional rendering-based pose estimation systems, such as those presented in (Wen et al., 2024) or (Labbé et al., 2022), is that our pose estimates indicate high uncertainty for distant, small, occluded, or non-visible objects as well dimensions along which the object is symmetric. A visualization of the pose beliefs at different points in the multi-step plan can be seen in Figure 14 in the Appendix.

A.5. Hyperparameters

Below we list out the significant hyperparameters involved in each baseline method, and how we chose them based on our hyperparameter search. We run the same seed for each hyperparameter and pick the best performing hyperparameter as the representative for our larger quantitative experiments in figure 3. The full domain randomization (FullDR) and no domain randomization (NoDR) baselines have no hyperparameters.

A.5.1. GoFLOW

We search over the following values of the α hyperparameter: [0.1, 0.5, 1.0, 1.5, 2.0]. We search over the following values of the β hyperparameters [0.0, 0.1, 0.5, 1.0, 2.0]. Other hyperparameters include number of network updates per training epoch ($K = 100$), network learning rate ($\eta_\phi = 1e - 3$), and neural spline flow architecture hyperparameters such as network depth ($\ell = 3$), hidden features (64), and number of bins (8). We implement our flow using the Zuko normalizing flow library (Rozet et al., 2022).

A.5.2. LSDR

Similarly to GoFlow, we search over the following values of the α_L hyperparameter: [0.1, 0.5, 1.0, 1.5, 2.0]. Other hyperparameters include the number of updates per training epoch ($T=100$), and initial Gaussian parameters: $\mu = (\xi_{max} + \xi_{min})/2.0$ and $\Sigma = \text{diag}(\xi_{max} - \xi_{min}/10)$

Method	FullDR	NoDR	DORAEMON	LSDR	ADR	GoFlow
Success Rate	6/10	3/10	5/10	5/10	5/10	9/10

Table 1. Real-world experimental results with the statistically significant results bolded.

	Cartpole	Ant	Quadcopter	Quadruped	Humanoid	Gears
FullDR	0.060±0.021	0.014±0.006	0.000±0.000	0.012±0.007	0.195±0.044	0.000±0.000
NoDR	0.005±0.001	0.000±0.000	0.000±0.000	0.086±0.008	0.001±0.001	0.008±0.003
DORAEMON	0.015±0.005	0.000±0.000	0.000±0.000	0.029±0.008	0.000±0.000	0.012±0.004
LSDR	0.103±0.013	0.000±0.000	0.007±0.003	0.009±0.008	0.176±0.082	0.000±0.000
ADR	0.006±0.001	0.000±0.000	0.000±0.000	0.070±0.006	0.000±0.000	0.008±0.004
GOFLOW	0.138±0.005	0.056±0.033	0.010±0.006	0.275±0.019	0.274±0.076	0.035±0.019

Table 2. Mean and standard error (SDE) of the final reward value, with entries in bold indicating the best-performing method or methods not statistically significantly worse than the best (one-tailed z test, $\alpha = 0.05$).

A.5.3. DORAEMON

We search over the following values of the ϵ_D hyperparameter: [0.005, 0.01, 0.05, 0.1, 0.5]. After fixing the best ϵ_D for each environment, we additionally search over the success threshold α_D : [0.005, 0.01, 0.05, 0.1, 0.5].

A.5.4. ADR

In ADR, we fix the upper threshold to be the success threshold $t^+ = J_T$ as was done in the original paper and search over the lower bound threshold $t^- = [0.1t^+, 0.25t^+, 0.5t^+, 0.75t^+, 0.9t^+]$. The value used in the original paper was $0.5t_H$. Other hyperparameters include the expansion/contraction rate, which we interpret to be a fixed fraction of the domain interval, $\Delta = 0.1 * [\xi_{\max} - \xi_{\min}]$, and boundary sampling probability $p_b = 0.5$.

A.6. Coverage vs Range Experiments

We compare coverage vs. range scale in the ant domain. We adjust the parameter lower and upper bounds outlined in Appendix A.3 and see how the coverage responds to those changes during training. The parameter range is defined relative to a nominal midpoint m set to the original domain parameters: $[m - (m - \text{lower}) * \text{scale}, m + (\text{upper} - m) * \text{scale}]$. The results of our experiment are shown in Figure 12

A.7. Coverage vs. Threshold Experiments

In this experiment, we assess the sensitivity of coverage to the performance threshold J_t for a single training seed. We experiment with threshold multipliers between 0 and $1.2 \cdot J_t$. The results of this experiment can be seen in Figure 13. These results demonstrate that although coverage is highly

dependent on J_t , GoFlow largely dominates the coverage distribution for wide ranges of J_t .

A.8. Real-world experiments

In addition to simulated experiments, we compare GoFlow against baselines on a real-world gear insertion task. In particular, we tested insertion of the pink medium gear over 10 trials for each baseline. To test this, we had the robot perform 10 successive pick/inserts of the pink gear into the middle shaft of the gear plate. Instead of randomizing the pose of the gear, we elected to fix the initial pose of the gear and the systematically perturb the end-effector pose by a random ± 0.01 meter translational offset along the x dimension during the pick. We expect some additional grasp pose noise due to position error during grasp and object shift during grasp. This led to a randomized in-hand offset while running the trained insertion policy. Our results show that GoFlow can indeed more robustly generalize to real-world robot settings under pose uncertainty.

A.9. Statistical Tests

We performed a statistical analysis of the simulated results reported in Figure 3 and the real-world experiments in Table 1. For the simulated results, we recorded the final domain coverages across all seeds and performed pairwise t-tests between each method and the top-performing method. The final performance mean and standard deviation are reported in Table 2. Any methods that were not significantly different from the top performing method ($p < 0.05$) are bolded. This same method was used to test significance of the real-world results.

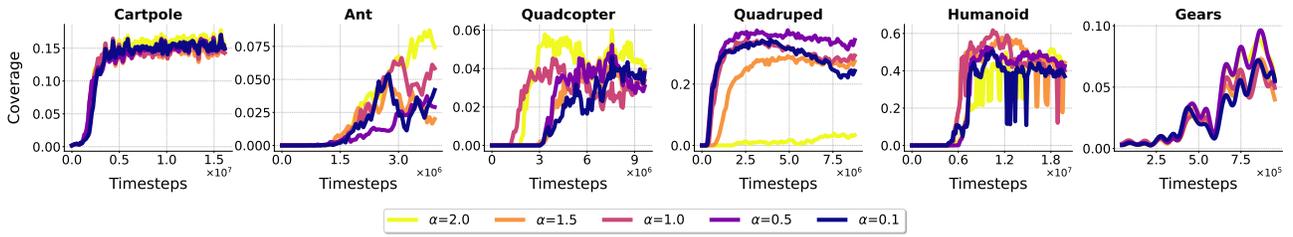


Figure 6. GoFlow hyperparameter sweep results for α

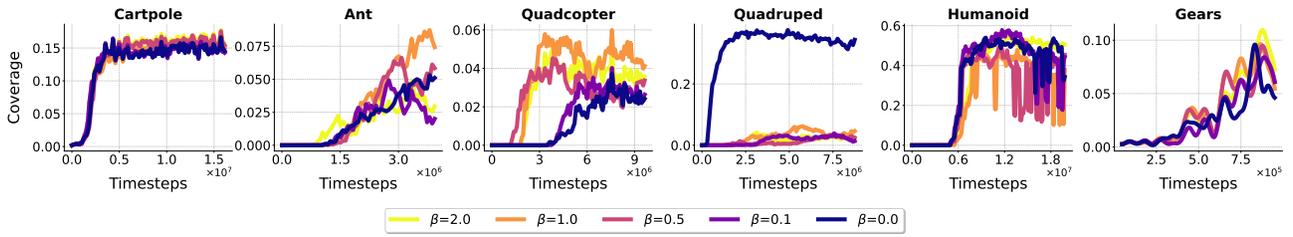


Figure 7. GoFlow hyperparameter sweep results for β after fixing the best α for each environment

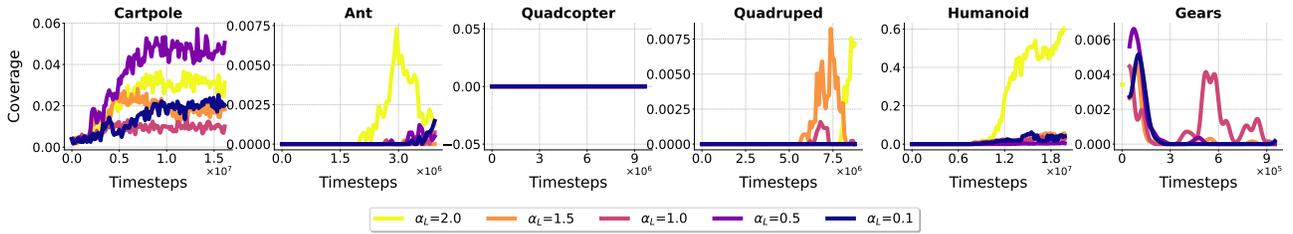


Figure 8. LSDR hyperparameter sweep results for α_L

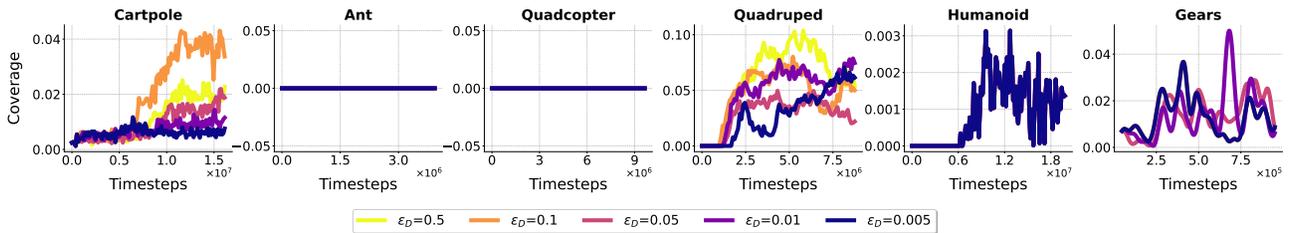


Figure 9. DORAEMON hyperparameter sweep results for ϵ_D

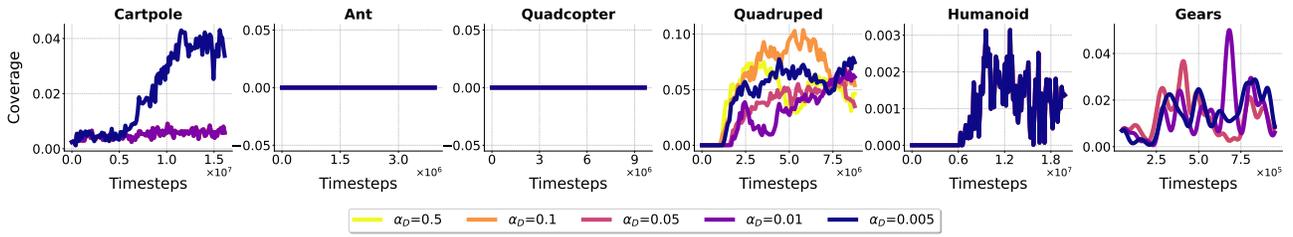


Figure 10. DORAEMON hyperparameter sweep results for α_D after fixing the best ϵ_D

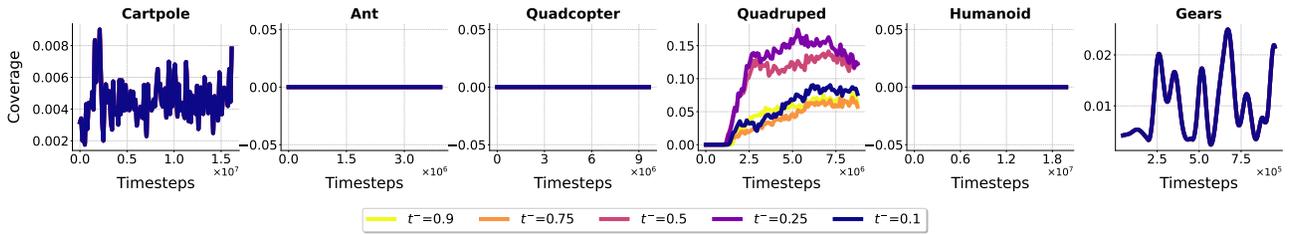


Figure 11. ADR hyperparameter sweep results for t^-

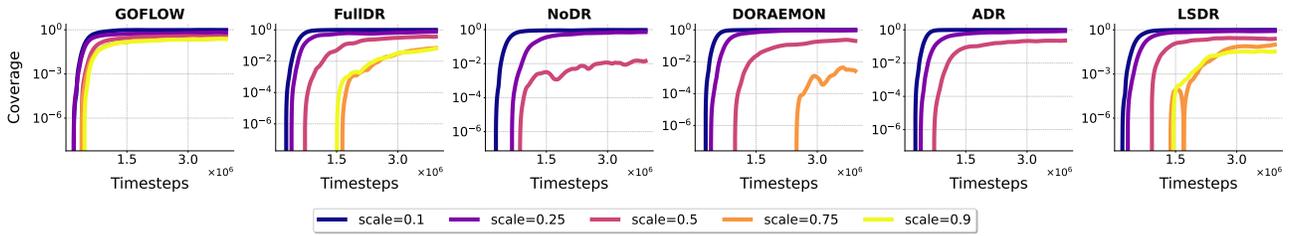


Figure 12. Coverage vs range experiment results discussed in Section A.6

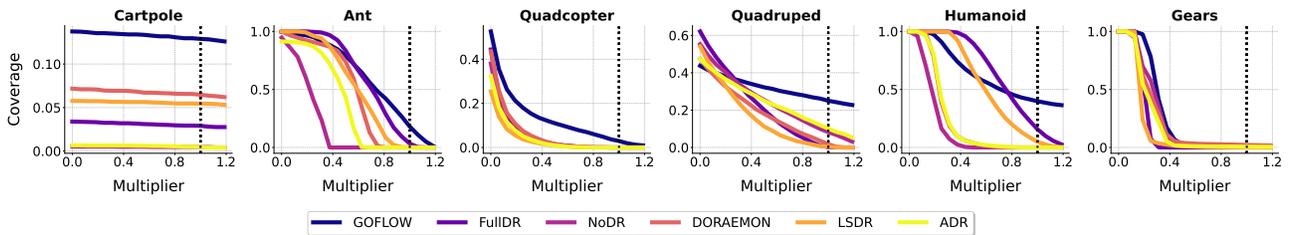


Figure 13. Coverage vs threshold experiment results discussed in Section A.7

	Cartpole	Ant	Quadcopter	Quadruped	Humanoid	Gears
FullDR	-920.132	362.454	-2.286	0.108	657.414	8.627
NoDR	-949.925	133.331	-3.807	-4.062	180.734	12.697
DORAEMON	-930.417	86.742	-1.767	-2.472	235.516	11.378
LSDR	-867.347	367.612	-1.870	-0.052	532.531	9.289
ADR	-953.498	64.815	-3.183	-2.926	246.297	9.671
GOFLOW	-820.510	198.302	-1.000	-2.814	376.652	14.209

Table 3. CVaR computed as the mean of the final rewards falling below the 10% percentile (VaR).

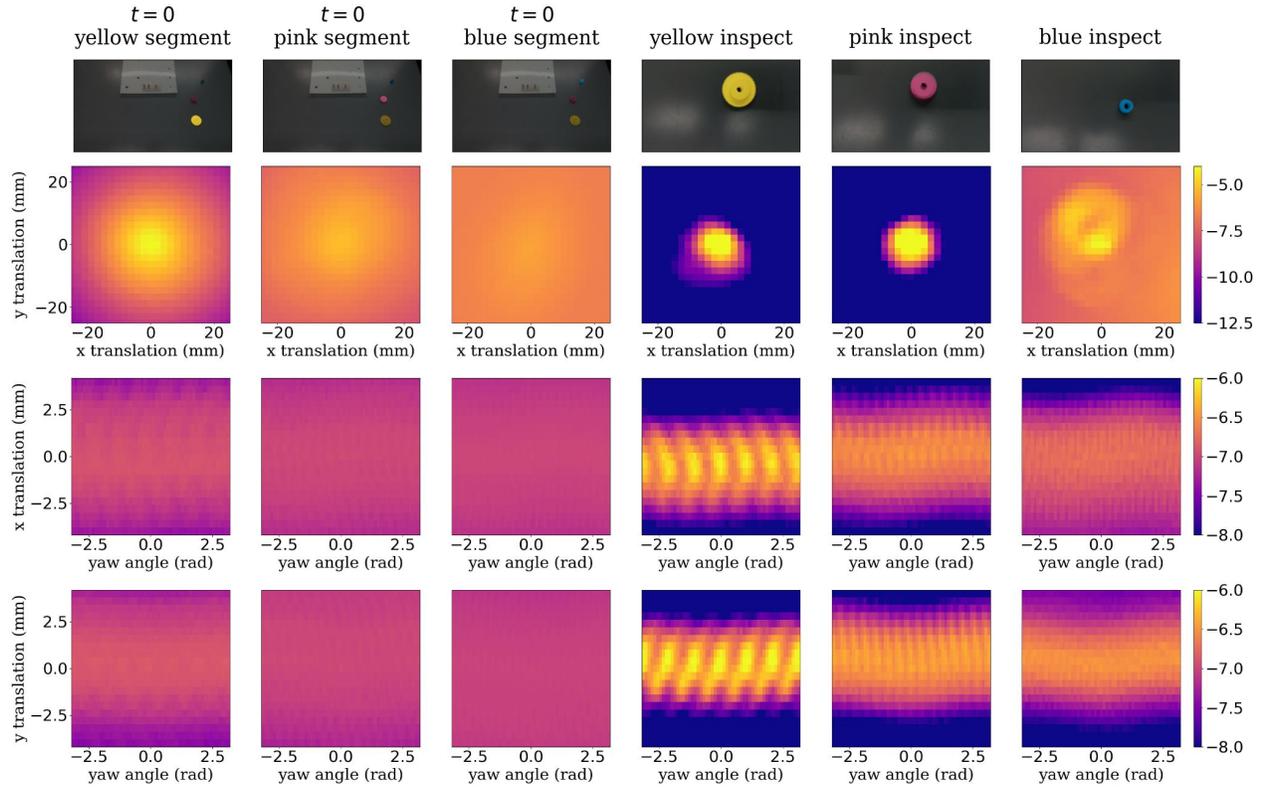


Figure 14. A visualization of the beliefs over the object pose under the initial image (first three columns) and after closer inspection (last three columns) as generated from the posterior of the model described in Section 6.3. The colormap corresponds to the log probability of the posterior pose estimate. All plots are centered around the most likely pose estimate under the image model.

Algorithm 2 Belief-Space Planner Using BFS

Require: Initial belief state b_0 , goal condition $G \subseteq \mathcal{B}$, set of skills \mathcal{A}_Π , success threshold η

```

1: Initialize the frontier  $\mathcal{F} \leftarrow \{b_0\}$ 
2: Initialize the visited set  $\mathcal{V} \leftarrow \emptyset$ 
3: Initialize the plan dictionary Plan mapping belief states to sequences of skills
4: while  $\mathcal{F}$  is not empty do
5:   Dequeue  $b$  from  $\mathcal{F}$ 
6:   if  $b \in G$  then
7:     return Plan[ $b$ ] ▷ Return the sequence of skills leading to  $b$ 
8:   end if
9:   for all skills  $\pi \in \mathcal{A}_\Pi$  do
10:    if  $b \in \text{Pre}_\pi$  given  $\eta$  then
11:       $b' \leftarrow \text{sample}(\text{Eff}_\pi)$ 
12:      if  $b' \notin \mathcal{V}$  then
13:        Add  $b'$  to  $\mathcal{F}$  and  $\mathcal{V}$ 
14:        Update Plan[ $b'$ ]  $\leftarrow$  Plan[ $b$ ] + [ $\pi$ ]
15:      end if
16:    end if
17:  end for
18: end while
19: return Failure ▷ No plan found

```
