

# Interactive Imitation Learning as Reinforcement Learning

*Perry Dong*



Electrical Engineering and Computer Sciences  
University of California, Berkeley

Technical Report No. UCB/EECS-2024-22

<http://www2.eecs.berkeley.edu/Pubs/TechRpts/2024/EECS-2024-22.html>

April 26, 2024

Copyright © 2024, by the author(s).  
All rights reserved.

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. To copy otherwise, to republish, to post on servers or to redistribute to lists, requires prior specific permission.

### Acknowledgement

I would like to thank my research advisors Professor Sergey Levine and Professor Yi Ma for their continued support and mentorship. The research this report is based on is part of a larger joint effort with Jianlan Luo, Yuexiang Zhai, Yi Ma, and Sergey Levine. I would like to thank all of them for their contributions.

*Masters Report*

---

**RLIF: Interactive Imitation Learning as Reinforcement Learning**

by Perry Dong

---

**Research Project**

Submitted to the Department of Electrical Engineering and Computer Sciences,  
University of California at Berkeley, in partial satisfaction of the requirements for the  
degree of **Master of Science, Plan II**.

Approval for the Report and Comprehensive Examination:

**Committee:**



---

Professor Sergey Levine  
Research Advisor

4/19/24

---

(Date)

\*\*\*\*\*



---

Professor Yi Ma  
Second Reader

4/23/24

---

(Date)

# Interactive Imitation Learning as Reinforcement Learning

Perry Dong

Spring 2024

# Abstract

Although reinforcement learning methods offer a powerful framework for automatic skill acquisition, for practical learning-based control problems in domains such as robotics, imitation learning often provides a more convenient and accessible alternative. In particular, an interactive imitation learning method such as DAgger, which queries a near-optimal expert to intervene online to collect correction data for addressing the distributional shift challenges that afflict naïve behavioral cloning, can enjoy good performance both in theory and practice without requiring manually specified reward functions and other components of full reinforcement learning methods. In this paper, we explore how off-policy reinforcement learning can enable improved performance under assumptions that are similar but potentially even more practical than those of interactive imitation learning. Our proposed method uses reinforcement learning with user intervention signals *themselves* as rewards. This relaxes the assumption that intervening experts in interactive imitation learning should be near-optimal and enables the algorithm to learn behaviors that improve over the potential suboptimal human expert. We then evaluate our method on challenging high-dimensional continuous control simulation benchmarks as well as real-world robotic vision-based manipulation tasks. The results show that it strongly outperforms DAgger-like approaches across the different tasks, especially when the intervening experts are suboptimal. Additional ablations also empirically show that the performance of our method is associated with the choice of intervention model and suboptimality of the expert.

## **Acknowledgements**

I would like to thank my research advisors Professor Sergey Levine and Professor Yi Ma for their continued support and mentorship. The research this report is based on is part of a larger joint effort with Jianlan Luo, Yuexiang Zhai, Yi Ma, and Sergey Levine. I would like to thank all of them for their contributions.

# Contents

<b>1</b>	<b>Introduction</b>	<b>4</b>
<b>2</b>	<b>Related Work</b>	<b>6</b>
<b>3</b>	<b>Background</b>	<b>7</b>
<b>4</b>	<b>Method</b>	<b>9</b>
<b>5</b>	<b>Experiments</b>	<b>11</b>
5.1	Intervention Strategies . . . . .	11
5.2	Performance on Continuous Control Benchmark Tasks . . . . .	12
5.3	Real-World Vision-Based Robotic Manipulation Task . . . . .	16
<b>6</b>	<b>Discussion</b>	<b>18</b>
<b>7</b>	<b>Acknowledgments</b>	<b>19</b>

# 1 Introduction

Reinforcement learning has achieved success in domains where well-specified reward functions are available, such as optimal control, games, and aligning large language models (LLMs) with human preferences [20, 15, 35, 29]. However, imitation learning methods are still often preferred

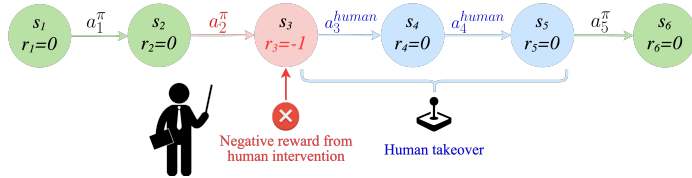


Figure 1: uses RL to learn without ground truth rewards, with data collected with suboptimal human interventions.

in some domains, such as robotics, because they are often more convenient, accessible, and easier to use. An often-cited weakness of naïve behavioral cloning is the compounding distributional shift induced by accumulating errors when deploying a learned policy. Interactive imitation learning methods, like the DAgger family of algorithms [34, 16, 13, 25, 33], address this issue by querying expert actions online and retraining the model iteratively in a supervised learning fashion. This performs well in practice, and in theory reduces the quadratic regret of imitation learning methods to be linear in the episode horizon. One particularly practical instantiation of this idea involves a human expert observing a learned policy, and intervening to provide *corrections* (short demonstrations) when the policy exhibits undesirable behavior [16, 37]. However, such interactive imitation learning methods still rely on interventions that are near-optimal, and offer no means to improve over the performance of the expert. Real human demonstrators are rarely optimal, and in domains such as robotics, teleoperation often does not afford the same degree of grace and dexterity as a highly tuned optimal controller. Can we combine the best parts of reinforcement learning and interactive imitation learning, combining the reward-maximizing behavior of RL, which can improve over the best available human behavior, with the accessible assumptions of interactive imitation learning?

The key insight we leverage in this work is that the decision to intervene during an interactive imitation episode itself can provide a reward signal for reinforcement learning, allowing us to instantiate RL methods that operate under similar but potentially weaker assumptions as interactive imitation methods, learning from human interventions but not assuming that such interventions are optimal. Intuitively, for many problems, it’s easier to detect a mistake than it is to optimally correct it. Imagine



an autonomous driving scenario with a safety driver. While the driver could intervene when the car deviates from good driving behavior, such interventions themselves might often be relatively uninformative and suboptimal – for example, if the human driver intervenes right before a collision by slamming on the breaks, simply teaching the policy to slam on the breaks is probably not the best solution, as it would be much better for the policy to learn to avoid the situations that necessitated such an intervention in the first place.

Motivated by this observation, we propose a method that runs RL on data collected from DAgger-style interventions, where a human operator observes the policy’s behavior and intervenes with *suboptimal* corrections when the policy deviates from optimal behavior. Our method labels the action that leads to an intervention with a negative reward and then uses RL to minimize the occurrence of intervention by maximizing these reward signals. We call our method *RLIF*: Reinforcement Learning via Intervention Feedback. This offers a convenient mechanism to utilize non-expert interventions: the final performance of the policy would not be bottlenecked by the suboptimality of the intervening expert, but rather the policy would improve to more optimally avoid interventions happening at all. Of course, the particular intervention strategy influences the behavior of such a method, and we require some additional assumptions on when interventions occur. We formalize several such assumptions and evaluate their effect on performance, finding that several reasonable strategies for selecting *when* to intervene lead to good performance.

Our main contribution is a practical RL algorithm that can be used under assumptions that closely resemble interactive imitation learning, without requiring ground truth reward signals. We empirically evaluate our approach in comparison to DAgger on a variety of challenging continuous control tasks, such as the Adroit dexterous manipulation and Gym locomotion environments [8]. Our empirical results show that our method is on average **2-3x** better than best-performing DAgger variants, and this difference is much more pronounced as the suboptimality gap expands. We also demonstrate our method scales to a challenging real-world robotic task involving an actual human providing feedback.

## 2 Related Work

**Interactive imitation learning.** Imitation learning extracts policies from static offline datasets via supervised learning [5, 2, 11, 28, 18]. Deploying such policies incurs distributional shift, because the states seen at deployment-time differ from those seen in training when the learned policy doesn’t perfectly match the expert, potentially leading to poor results [32, 34]. Interactive imitation learning leverages additional online human interventions from states visited by the learned policy to address this issue [13, 16, 12, 25, 34]. These methods generally assume that the expert interventions are near-optimal. Our method relaxes this assumption, by using RL to train on data collected in this interactive fashion, with rewards derived from the user’s choice of when to intervene.

**Imitation learning with reinforcement learning.** Another line of related work uses RL to improve on suboptimal human demonstrations [40, 38, 6, 39, 26, 1, 30, 17, 23, 41, 42, 36, 3]. These methods typically initialize the RL replay buffer with human demonstrations, and then improve upon those human demonstrations by running RL with the task reward. In contrast to these methods, our approach does not require any task reward, but rather recovers a reward signal implicitly from intervention feedback. Some works use RL with interventions but assume the expert is optimal [22], which our method does not assume. Other works incorporate example high-reward states specified by a human user in place of demonstrations [31, 7]. While this is related to our approach of assigning negative rewards at intervention states, our interventions are collected interactively during execution under assumptions that match interactive imitation learning, rather than being provided up-front. Closely related to our approach, Kahn et al. [14] proposed a robotic navigation system that incorporates *disengagement* feedback, where a model is trained to predict states where a user will halt the robot, and then avoids those states. Our framework is model-free and operates under general interactive imitation assumptions, and utilizes more standard DAgger-style interventions rather than just disengagement signals.

### 3 Background

In this section, we set up the relevant background for interactive imitation learning and reinforcement learning and then introduce our problem statement.

**Behavioral cloning and interactive imitation learning.** The most basic form of imitation learning is behavioral cloning, which simply trains a policy  $(a|s)$  on a dataset of demonstrations  $D$ , conventionally assumed to be produced by an optimal expert policy  $(a|s)$ , with  $d(s)$  as its state marginal distribution. Then, for each  $(s, a) \in D$ , behavioral cloning assumes that  $s \sim d(s)$  and  $a \sim (a|s)$ . Behavioral cloning then chooses  $\pi = \arg \min_{\pi \in \Pi} \sum_{s, a \in D} \ell(s, a, \pi)$ , where  $\ell(s, a, \pi)$  is some loss function, such as the negative log-likelihood (i.e.,  $\ell(s, a, \pi) = -\log \pi(a|s)$ ). Naïve behavioral cloning is known to accumulate regret quadratically in the time horizon  $H$ : when  $\pi$  differs from  $\pi^*$  even by a small amount, erroneous actions will lead to distributional shift in the visited states, which in turn will lead to larger errors [34]. Interactive imitation learning methods, such as DAgger and its

variants [34, 32, 16, 13, 25, 12], propose to address this problem, reducing the error to be *linear* in the horizon by gathering additional training data by running the learned policy  $\pi$ , essentially adding new samples  $(s, a)$  where  $s \sim d(s)$ , and  $a \sim (a|s)$ . Different interactive imitation learning methods prescribe different strategies for adding such labels. Classic DAgger [34] runs  $\pi$  and then asks a human expert to relabel the resulting states with  $a \sim (a|s)$ . This is often unnatural in time-sensitive control settings, such as

robotics and driving, and a more user-friendly alternative such as HG-DAgger and its variants [16] instead allows a human expert to *intervene*, taking over control from  $\pi$  and overriding it with an expert action. We illustrate this in Algorithm 1.

Although this changes the state distribution, the essential idea of the method (and its regret bound) remain the same. Our aim in this work will be to apply RL to this setting to address this issue, potentially even outperforming the expert.

---

**Algorithm 1** Interactive imitation

---

**Require:**  $\pi, \pi^*, D$

```
1: for trial  $i = 1$  to  $N$  do
2:   Train  $\pi$  on  $D$  via supervised learning
3:   for timestep  $t = 1$  to  $T$  do
4:     if  $\pi^*$  intervenes at  $t$  then
5:       append  $(s_t, a_t)$  to  $D_i$ 
6:     end if
7:   end for
8:    $D \leftarrow D \cup D_i$ 
9: end for
```

---

**Reinforcement learning.** RL algorithms aim to learn optimal policies in Markov decision processes (MDPs). We will use an infinite-horizon formulation in our analysis. The MDP is defined as  $M = \{S, A, P, r, \gamma\}$ .  $M$  comprises:  $S$ , a state space of cardinality  $S$ ,  $A$ , an action space with size  $A$ ,  $P : S \times A \rightarrow \Delta(S)$ , representing the transition probability of the MDP,  $r : S \times A \rightarrow [0, ]$  is the reward function, and  $\gamma \in (0, 1)$  represents the discount factor. We use  $\pi : S \rightarrow \Delta(A)$ . We introduce the value function  $V^\pi(s)$  and the Q-function  $Q^\pi(s, a)$  associated with policy  $\pi$  as:  $V^\pi(s) := E \sum_{t=0}^{\infty} \gamma^t r(s_t, a_t) | s_0 = s; \pi, \forall s \in S$  and  $\forall (s, a) \in S \times A : Q^\pi(s, a) := E \sum_{t=0}^{\infty} \gamma^t r(s_t, a_t) | s_0 = s, a_0 = a; \pi$ , as is standard in reinforcement learning analysis. We assume the initial state distribution is given by  $\mu: s_0 \sim \mu$ , and  $\mu \in \Delta(S)$  and we slightly abuse the notation by using  $V^\pi(\mu)$  to denote  $E_{s \sim \mu} V^\pi(s)$ . The goal of RL is to learn an optimal policy in the policy class  $\Pi$  that maximizes the expected cumulative reward within the horizon  $H: = \arg \max_{\pi \in \Pi} V^\pi(\mu)$  [4]. Without loss of generality, we assume the optimal policy to be *deterministic* [4, 21]. We slightly abuse the notation by using  $V^*, Q^*$  to denote  $V, Q$ . Additionally, we use  $d_\mu^\pi(s) = (1 - \gamma) \sum_{t=0}^{\infty} \gamma^t P^\pi(s_t = s | s_0 \sim \mu)$ , to denote the state occupancy distribution under policy  $\pi$  on the initial state distribution  $s_0 \sim \mu$ . We also slightly abuse the notation by using  $d_\mu^\pi \in R^S$  to denote a vector, whose entries are  $d_\mu^\pi(s)$ .

**Problem setup.** Our aim will be to develop a reinforcement learning algorithm that operates under assumptions that resemble interactive imitation learning, where the algorithm is not provided with a reward function, but instead receives demonstrations followed by interactive interventions, as discussed above. We will not assume that the actions in the interventions themselves are optimal, but will make an additional mild assumption that the choice of *when* to intervene itself carries valuable information. We will discuss the specific assumption used in our analysis in Section ??, and we utilize several intervention strategies in our experiments, but intuitively we assume that the expert is more likely to intervene when takes a bad action. This in principle can provide an RL algorithm with a signal to alter its behavior, as it suggests that the steps leading up to this intervention deviated significantly from optimal behavior. Thus, we will aim to relax the strong assumption that the expert is optimal in exchange for the more mild and arguably natural assumption that the expert’s *choice of when to intervene* correlates with the suboptimality of the learned policy.

## 4 Method

The key observation of this paper is that, under typical interactive imitation learning settings, the *intervention signal alone* can provide useful information for RL to optimize against, without assuming that the expert demonstrator actually provides optimal actions.

Our method, which we refer to as RLIF (reinforcement learning from intervention feedback), follows a similar outline as Algorithm 1, but using reinforcement learning in place of supervised

learning. The generic form of our method is provided in Algorithm 2. On Line 2, the policy is trained on the aggregated dataset  $D$  using RL, with rewards derived from interventions. The reward is simply set to 0 for any transition where the expert does not intervene, and -1 for the previous transition where an expert intervenes. After an intervention, the expert can also optionally take over the control for a few steps; after which it will be released to the RL agent. This way, an expert can also intervene multiple times during one episode. All transitions are added to the dataset  $D$ , not just those that contain the expert reward; in practice, we can optionally initialize  $D$  with a small amount of offline data to warm-start the process. An off-policy RL algorithm can utilize all of the data and can make use of the reward labels to avoid situations that cause interventions. This approach has a number of benefits: unlike RL, it doesn't require the true task reward to be specified, and unlike interactive imitation learning, it does not require the expert interventions to contain optimal actions, though it does require the *choice* of when to intervene to correlate with the suboptimality of the policy (as we will discuss later). Intuitively, we expect it to be less of a burden for experts to *only* point out which states are undesirable rather than actually *act optimally* in those states.

---

**Algorithm 2** RLIF

---

**Require:**  $\pi, \mathcal{D}, D$ 

```
1: for trial  $i = 1$  to  $N$  do
2:   Train  $\pi$  on  $D$  via reinforcement learning.
3:   for timestep  $t = 1$  to  $T$  do
4:     if intervenes at  $t$  then
5:       label  $(s_{t-1}, a_{t-1}, s_t)$  with -1 reward,
        append to  $D_i$ 
6:     else
7:       label  $(s_{t-1}, a_{t-1}, s_t)$  with 0 reward,
        append to  $D_i$ 
8:     end if
9:   end for
10:   $D \leftarrow D \cup D_i$ 
11: end for
```

---

**Practical implementation.** In order to instantiate Algorithm 2 in a practical deep RL framework, the base RL algorithm must be chosen to take advantage of a combination of on-policy samples and potentially suboptimal near-expert interventions. This necessitates using a off-policy RL algorithm that can incorporate prior (near-expert) data easily but also can efficiently improve with online experience. While a variety of algorithms designed for online RL with offline data could be suitable [36, 19, 27], we adopt the recently proposed RLPD algorithm [3], which has shown compelling results on sample-efficient robotic learning. RLPD is an off-policy actor-critic reinforcement learning algorithm that builds on soft-actor critic [10], but makes some key modifications to satisfy the desiderata above such as a high update-to-data ratio, layer-norm regularization during training, and using ensembles of value functions, which make it more suitable for incorporating offline data into online RL. For further details on this method, we refer readers to prior work [3], though we emphasize that our method is generic and in principle could be implemented relatively easily on top of a variety of RL algorithms. The RL algorithm itself is not modified, and our method can be viewed as a meta-algorithm that simply changes the dataset on which the RL algorithm operates.

## 5 Experiments

Since our method operates under standard interactive imitation learning assumptions, our experiments aim to compare RLIF to DAgger under different types of suboptimal experts and intervention modes. We seek to answer the following questions: (1) Are the intervention rewards sufficient signal for RL to learn effective policies? (2) How well does our method perform compared to DAgger, especially with suboptimal experts? (3) What are the implications of different intervention strategies on empirical performance?

### 5.1 Intervention Strategies

In order to learn from interventions, we need the intervening experts to convey useful information about the task through their decision about *when* to intervene. Since real human experts are likely to be imperfect in making this decision, we study a variety of intervention strategies in simulation empirically to validate the stability of our method with respect to this assumption. The baseline strategy, which we call **Random Intervention**, simply intervenes uniformly at random, with equal probability at every step. The more intelligent model, **Value-Based Intervention**, strategy assumes the expert intervenes with probability  $\beta$  when there is a gap  $\delta$  between actions from the expert and agent w.r.t. a reference value function  $Q$ . This model aims to capture an uncertain and stochastic expert, who might have a particular policy that they use to choose actions (which could be highly suboptimal), and a *separate* value function  $Q$  with its corresponding policy that they use to determine if the robot is doing well or not, which they use to determine when to intervene. Note that might be much better than – for example, a human expert might correctly determine the robot is choosing poor actions for the task, even if their own policy is not good enough to perform the task either.  $\delta$  represents the confidence level of the intervening expert: the smaller  $\delta$  is, the more willing they are to intervene on even slightly suboptimal robot actions. We formalize

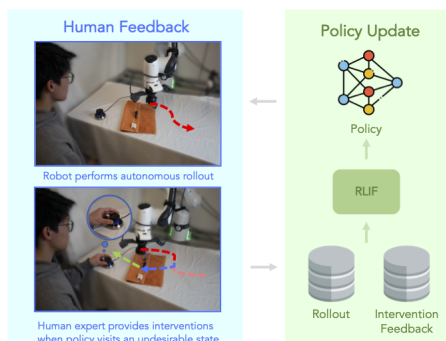


Figure 2: A human operator supervises policy training and provides intervention with a 3D mouse.

this model in Eq. 1. In practice, we choose a value for  $\beta$  close to 1, such as 0.95.

$$P(\text{Intervention}|s) = \begin{cases} \beta, & \text{if } Q(s, (s)) > Q(s, \pi(s)) + \delta \\ 1 - \beta, & \text{otherwise.} \end{cases} \quad (1)$$

This model may not be fully representative of real human behavior, so we also evaluate our method with real human interventions, where an expert user provides intervention feedback to a real-world robotic system performing a peg insertion task, shown in Figure 2. We discuss this further in Sec. 5.3.

## 5.2 Performance on Continuous Control Benchmark Tasks

First, we evaluate RLIF in comparison with various interactive imitation learning methods on several high-dimensional continuous control benchmark tasks. These experiments vary both the optimality of the expert’s policies and the expert’s intervention strategy.

**Simulation experiment setup.** We use Gym locomotion and Adroit dexterous manipulation tasks in these experiments, based on the D4RL environments [9]. The Adroit environments require controlling a 24-DoF robotic hand to perform tasks such as opening a door or rotating a pen to randomly sampled goal locations. The Gym locomotion task (walker) requires controlling a planar 2-legged robot to walk forward. Both domains require continuous control at relatively high frequencies, where approximation errors and distributional shift can accumulate quickly for naïve imitation learning methods. Note that although these benchmarks are often used to evaluate the performance of RL algorithms, we do not assume access to any reward function beyond the signal obtained from the expert’s interventions, and therefore our main point of comparison are DAgger variants, rather than prior RL algorithms.

The experts and  $Q$  associated with intervention strategies are obtained by training

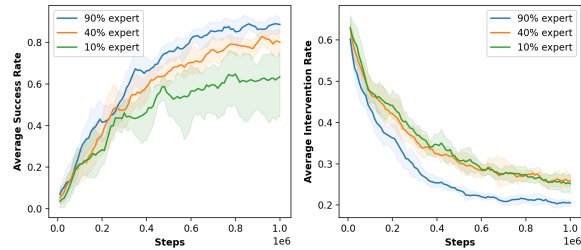


Figure 3: Average success rate and intervention rate for the Adroit-Pen task during training, as the agent improves, the intervention decreases.



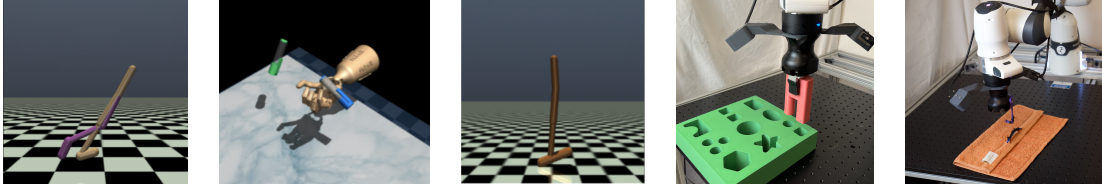


Figure 4: **Tasks in our experimental evaluation:** Benchmark tasks Walker2d, Pen, and Hopper and two vision-based contact-rich manipulation tasks on a real robot. The benchmark tasks require handling complex high-dimensional dynamics and underactuation. The robotic insertion task requires additionally addressing complex inputs such as images, non-differentiable dynamics such as contact, and all sensor noise associated with real-world robotic settings.

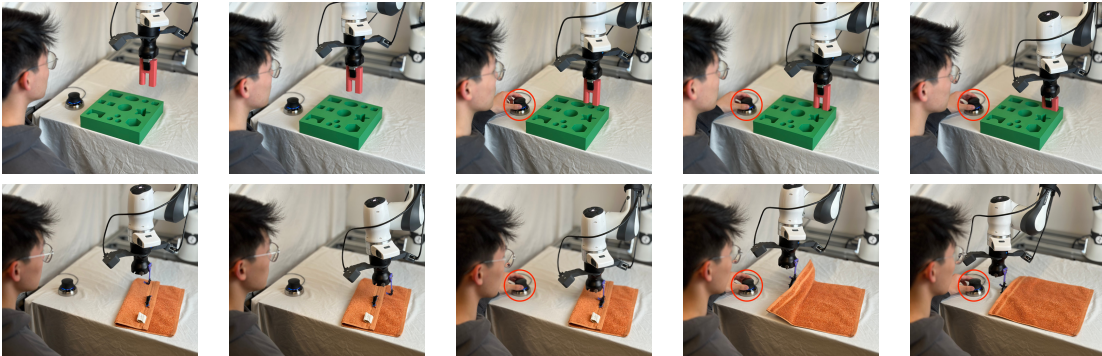


Figure 5: Sequential steps of robot manipulation for the peg insertion and cloth unfolding tasks on a real robot.

policies on subsampled datasets to induce the desired level of suboptimality. To perform controlled experimentation of RLIF against DAgger and HG-DAgger, we prepare offline datasets to use for pretraining on DAgger and HG-DAgger and to initialize to replay buffer for RLIF. The initial datasets of all simulation tasks are subsets of datasets provided in d4rl. The specific dataset used to subsample and sizes of the initial dataset for each task are listed in Table 2. The rewards for all timesteps for all of the initial datasets are set to zero.

Domain	Expert Level	RLIF with Value Based Intervention	RLIF with Random Intervention	HG-Dagger	HG-Dagger with 85% Random Intervention	Dagger	Dagger with 85% Random Intervention	BC
adroit-pen	~90%	<b>88.47±3.06</b>	42.87±12.86	73.47±6.19	74.27±5.79	78.13±3.24	79.07±9	54.13±14.24
	~40%	<b>80.87±6.01</b>	34.13±10.32	60±3.58	29.33±6.56	35.73±7.49	38.67±2.06	
	~10%	<b>64.04±17.59</b>	28.33±4.43	28.53±7.66	9.47±4.09	8.93±1.43	12.8±6.25	
	average	<b>77.79±8.89</b>	35.11±9.20	54±5.81	37.69±5.48	40.93±4.05	43.51±5.77	
locomotion-walker2d	~110%	108.99±5.28	106.51±0.47	53.55±9.76	<b>112.7±2.51</b>	57.94±8.69	76.13±3.27	44.46±13.59
	~70%	<b>99.66±5.9</b>	75.62±50.02	44.75±3.46	69.73±5.99	20.49±3.15	43.59±2.56	
	~20%	<b>102.85±2.26</b>	19.11±24.08	11.94±0.88	19.66±3.69	12.37±2.96	20.1±2.17	
	average	<b>103.83±4.48</b>	67.08±43.83	36.75±4.7	67.36±4.06	30.27±4.93	46.61±2.67	
locomotion-hopper	~110%	<b>109.17±0.16</b>	93.76±7.87	80.3±14.74	86.93±4.85	70.58±9.98	61.64±11.36	64.77±10.23
	~40%	<b>108.42±0.62</b>	103.9±10.28	40.66±3.35	42.65±2.86	38.7±3.7	19.63±2.3	
	~15%	<b>108.01±0.64</b>	75.12±28.95	25.2±3.58	24.37±2.26	19.54±2.14	10.29±1.24	
	average	<b>108.53±0.47</b>	90.93±15.7	48.72±7.22	51.32±3.32	42.94±5.27	30.46±4.97	

Table 1: A comparison of RLIF, HG-Dagger, Dagger, and BC on continuous control tasks. RLIF consistently performs better than HG-Dagger and Dagger baselines for each individual expert level as well as averaged over all expert levels.

Tasks	Dataset	Subsampled Size
Adroit Pen	pen-expert-v1	50 trajectories
Locomotion Hopper	hopper-expert-v2	50 trajectories
Locomotion Walker2d	walker2d-expert-v2	10 trajectories
Real-World Robot Experiments	manually collected suboptimal trajectories	5 trajectories

Table 2: Offline datasets for each task.

Experts of varying levels are trained for all tasks on the dataset with either BC, IQL, SAC, or RLPD depending on the task. The various levels of the experts are obtained by training on the human dataset or expert datasets subsampled to various sizes. For each level and task, the same expert is used to intervene across all intervention strategies for both RLIF, Dagger, and HG-Dagger.

**Results and discussion.** We report our main results in Table. 1. Each table depicts different expert performance levels on different tasks (rows), and different methods with different intervention strategies (columns). Note the “expert level” in the table

indicates the suboptimality of the expert; for example, a “90%” expert means it can achieve 90% of the reference optimal expert score of a particular task. This score is normalized across environments to be 100.0. Since we trained such experts using our curated datasets, the normalized score can exceed 100.0. For the rest of the table, we report the performance of different methods w.r.t. this normalized reference score.

To start the learning process, we initialize the replay buffer with a small number of samples from the dataset curation method described above. We compare RLIF with DAgger and HG-DAgger under different intervention modes and expert levels. Specifically, we run DAgger under two interven-

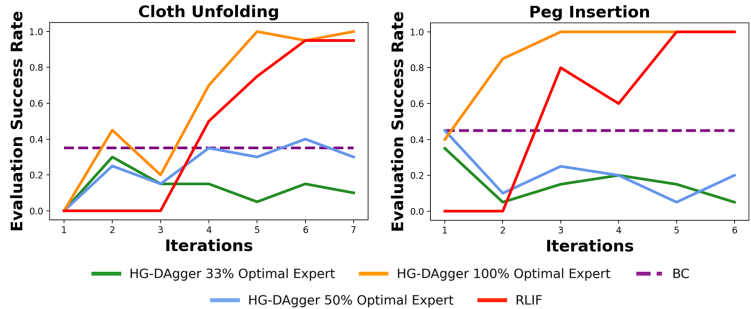


Figure 6: RLIF on the real-world robotic manipulation task.

tion modes: 1) the expert issues an intervention if the difference between the agent’s action and the expert’s action is larger than a threshold, as used in [16, 34], 2) the expert issues random intervention uniformly with a given probability at each step of the episode. The results suggest several takeaways: (1) We see that regardless of the suboptimality of the experts, RLIF with value-based interventions can indeed reach good performance, even if the suboptimality gap is *very large*, such as improving a 15% expert to a score over 100.0 in the 2D-walker task. (2) RLIF with a value-based simulated interventions outperforms RLIF with random interventions, especially when the suboptimality is large. This confirms that the interventions carry a meaningful signal about the task, providing implicit rewards to RLIF. We also observe that the value-based intervention rate goes down as the agent improves, as shown in Fig. 3. This confirms our proposed intervention model does carry useful information about the task, and works reasonably as the agent learns. (3) RLIF outperforms both DAgger and HG-DAgger consistently across all tasks, and the performance gap is more pronounced when the expert is *more suboptimal*, and can be as large as **5x**; crucially, this resonates with our motivation: the performance of DAgger-like algorithms will be subject to the suboptimality of the experts, while our method can reach good performance even with suboptimal experts by learning from the expert’s decision of *when* to intervene.

**Ablations on Intervention Modes.** As stated in Sec. 5.1, the performance of RLIF critically depends on when the expert chooses to intervene. Now we analyze the effect of different  $Q$  value functions on the final performance of RLIF. We report the numbers in Table. 3. We can observe that the particular choice of  $Q$  does heavily influence our algorithm’s performance: as the deteriorates,  $Q$  becomes increasingly inaccurate, making the intervention decision more “uncalibrated.” This translates to worse policy performance.

Domain	Expert Level	10% $Q$	60% $Q$	110% $Q$
locomotion	~110%	78.9±49.97	85.91±19.31	108.99±5.28
-walker	~40%	38.78±20.3	71.21±21.96	99.66±5.9
	~20%	33.55±9.63	66.8±8.22	102.85±2.26

Table 3: An ablation of RLIF on walker2d with different  $Q$  and expert levels. A 10%  $Q$  means that it was trained with the offline dataset generated by a 10% expert.

### 5.3 Real-World Vision-Based Robotic Manipulation Task

While we have shown that RLIF works well under reasonable models of the expert’s intervention strategy, to confirm that this model actually reflects real human interventions, we next conduct an experiment where a human operator supplies the interventions directly for a real-world robotic manipulation task. The first task, shown in Fig. 2, involves controlling a 7-DoF robot arm to fit a peg into its matching shape with a very tight tolerance (1.5mm), directly from image observations. This task is difficult because it involves dealing with discontinuous and non-differentiable contact dynamics, complex high-dimensional inputs from the camera, and an imperfect human operator. The second task involves controlling the same robot arm to unfold a piece of cloth by hooking onto it and pulling it open. This task is difficult because it requires manipulation of a deformable object: specifying rewards programmatically for such a task is very challenging, and the perception system trained via RL must be able to keep track of the complex object geometry. Filmstrips of the tasks are shown in Fig. 5 to visualize the sequential steps of both the peg insertion and cloth unfolding tasks.

For the insertion task, a trial is counted as a success if the peg is inserted into its matching hole with a certain tolerance, and for the unfolding task, a trial is counted as a success if the cloth is successfully unfolded all the way. All success rates are reported based on 20 trials. We report the results in Fig. 6. Our method can solve the

insertion task with a 100% success rate within six rounds of interactions. It solves the unfolding task with a 95% success rate in seven rounds of interaction. This highlights the practical usability of our method in challenging real-world robotic tasks.

## 6 Discussion

We present a reinforcement learning method that learns from interventions in a setting that closely resembles interactive imitation learning without assumptions of expert optimality and access to ground truth rewards. However, our approach does have a number of limitations. First, we require an RL method that can actually train on all of the data collected by both the policy and the intervening expert. We were able to use an off-the-shelf offline RL algorithm (RLPD) without modification, but in general this is not an easy reinforcement learning problem. Second, imitation learning is often preferred precisely because it doesn't require online deployment. While in practice deploying a policy under expert oversight might still be safer (e.g., with a safety driver in the case of autonomous driving), investigating the safety of online deployment of RL algorithms particularly in an interactive learning setting is an important direction for future work.

## 7 Acknowledgments

Thanks to Jianlan Luo, Yuexiang Zhai, Yi Ma, and Sergey Levine for their collaborations on this project. This report is a select subset of work that was originally published as the conference paper “RLIF: Interactive Imitation Learning as Reinforcement Learning” [24] at the 2024 International Conference on Learning Representations (ICLR).

## References

- [1] Samuel Ainsworth, Matt Barnes, and Siddhartha Srinivasa. Mo’s states mo’problems: Emergency stop mechanisms from observation. *Advances in Neural Information Processing Systems*, 32, 2019.
- [2] Brenna D. Argall, Sonia Chernova, Manuela Veloso, and Brett Browning. A survey of robot learning from demonstration. *Robot. Auton. Syst.*, 57(5):469–483, may 2009. ISSN 0921-8890. doi: 10.1016/j.robot.2008.10.024. URL <https://doi.org/10.1016/j.robot.2008.10.024>.
- [3] Philip J. Ball, Laura Smith, Ilya Kostrikov, and Sergey Levine. Efficient online reinforcement learning with offline data. In Andreas Krause, Emma Brunskill, Kyunghyun Cho, Barbara Engelhardt, Sivan Sabato, and Jonathan Scarlett, editors, *Proceedings of the 40th International Conference on Machine Learning*, volume 202 of *Proceedings of Machine Learning Research*, pages 1577–1594. PMLR, 23–29 Jul 2023. URL <https://proceedings.mlr.press/v202/ball123a.html>.
- [4] Dimitri Bertsekas. *Reinforcement learning and optimal control*. Athena Scientific, 2019.
- [5] Aude Billard, Sylvain Calinon, Rüdiger Dillmann, and Stefan Schaal. *Robot Programming by Demonstration*, pages 1371–1394. Springer Berlin Heidelberg, Berlin, Heidelberg, 2008. ISBN 978-3-540-30301-5. doi: 10.1007/978-3-540-30301-5\_60. URL [https://doi.org/10.1007/978-3-540-30301-5\\_60](https://doi.org/10.1007/978-3-540-30301-5_60).
- [6] Ching-An Cheng, Xinyan Yan, Nolan Wagener, and Byron Boots. Fast policy learning through imitation and reinforcement. 2018.
- [7] Ben Eysenbach, Sergey Levine, and Russ R Salakhutdinov. Replacing rewards with examples: Example-based policy search via recursive classification. *Advances in Neural Information Processing Systems*, 34:11541–11552, 2021.
- [8] J. Fu, A. Kumar, O. Nachum, G. Tucker, and S. Levine. D4rl: Datasets for deep data-driven reinforcement learning. In *arXiv*, 2020. URL <https://arxiv.org/pdf/2004.07219>.



- [9] Justin Fu, Aviral Kumar, Ofir Nachum, George Tucker, and Sergey Levine. D4rl: Datasets for deep data-driven reinforcement learning. *arXiv preprint arXiv:2004.07219*, 2020.
- [10] T. Haarnoja, A. Zhou, P. Abbeel, and S. Levine. Soft actor-critic: Off-policy maximum entropy deep reinforcement learning with a stochastic actor. In *arXiv*, 2018. URL <https://arxiv.org/pdf/1801.01290.pdf>.
- [11] J. Ho and S. Ermon. Generative adversarial imitation learning. In *Neural Information Processing Systems (NIPS)*, 2016.
- [12] Ryan Hoque, Ashwin Balakrishna, Ellen R. Novoseller, Albert Wilcox, Daniel S. Brown, and Ken Goldberg. Thriftydagger: Budget-aware novelty and risk gating for interactive imitation learning. In *Conference on Robot Learning*, 2021.
- [13] Ryan Hoque, Lawrence Yunliang Chen, Satvik Sharma, K Dharmarajan, Brijen Thananjeyan, P. Abbeel, and Ken Goldberg. Fleet-dagger: Interactive robot fleet learning with scalable human supervision. In *Conference on Robot Learning*, 2022.
- [14] Gregory Kahn, P. Abbeel, and Sergey Levine. Land: Learning to navigate from disengagements. *IEEE Robotics and Automation Letters*, 6:1872–1879, 2020.
- [15] Dmitry Kalashnikov, Alex Irpan, Peter Pastor, Julian Ibarz, Alexander Herzog, Eric Jang, Deirdre Quillen, Ethan Holly, Mrinal Kalakrishnan, Vincent Vanhoucke, et al. Scalable deep reinforcement learning for vision-based robotic manipulation. In *Conference on Robot Learning*, pages 651–673, 2018.
- [16] Michael Kelly, Chelsea Sidrane, K. Driggs-Campbell, and Mykel J. Kochenderfer. Hg-dagger: Interactive imitation learning with human experts. *2019 International Conference on Robotics and Automation (ICRA)*, pages 8077–8083, 2018.
- [17] Rahul Kidambi, Aravind Rajeswaran, Praneeth Netrapalli, and Thorsten Joachims. Morel: Model-based offline reinforcement learning. *arXiv preprint arXiv:2005.05951*, 2020.
- [18] Michael Laskey, Jonathan Lee, Roy Fox, Anca Dragan, and Ken Goldberg. Dart: Noise injection for robust imitation learning, 2017.

- [19] Seunghyun Lee, Younggyo Seo, Kimin Lee, Pieter Abbeel, and Jinwoo Shin. Offline-to-online reinforcement learning via balanced replay and pessimistic q-ensemble. In *Conference on Robot Learning*, pages 1702–1712. PMLR, 2022.
- [20] Sergey Levine, Chelsea Finn, Trevor Darrell, and Pieter Abbeel. End-to-end training of deep visuomotor policies. *The Journal of Machine Learning Research*, 17(1):1334–1373, 2016.
- [21] Gen Li, Laixi Shi, Yuxin Chen, Yuejie Chi, and Yuting Wei. Settling the sample complexity of model-based offline reinforcement learning. *arXiv preprint arXiv:2204.05275*, 2022.
- [22] Quanyi Li, Zhenghao Peng, and Bolei Zhou. Efficient learning of safe driving policy via human-ai copilot optimization, 2022.
- [23] Jianlan Luo, Oleg O. Sushkov, Rugile Pevceviciute, Wenzhao Lian, Chang Su, Mel Vecerík, Ning Ye, Stefan Schaal, and Jonathan Scholz. Robust multi-modal policies for industrial assembly via reinforcement learning and demonstrations: A large-scale study. *ArXiv*, abs/2103.11512, 2021.
- [24] Jianlan Luo, Perry Dong, Yuexiang Zhai, Yi Ma, and Sergey Levine. Rlif: Interactive imitation learning as reinforcement learning, 2024.
- [25] Kunal Menda, K. Driggs-Campbell, and Mykel J. Kochenderfer. Ensembledagger: A bayesian approach to safe imitation learning. *2019 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, pages 5041–5048, 2018.
- [26] Ashvin Nair, Bob McGrew, Marcin Andrychowicz, Wojciech Zaremba, and Pieter Abbeel. Overcoming exploration in reinforcement learning with demonstrations. In *2018 IEEE International Conference on Robotics and Automation (ICRA)*, pages 6292–6299, 2018. doi: 10.1109/ICRA.2018.8463162.
- [27] Mitsuhiko Nakamoto, Yuexiang Zhai, Anikait Singh, Max Sobol Mark, Yi Ma, Chelsea Finn, Aviral Kumar, and Sergey Levine. Cal-ql: Calibrated offline rl pre-training for efficient online fine-tuning. *arXiv preprint arXiv:2303.05479*, 2023.
- [28] Takayuki Osa, Joni Pajarinen, Gerhard Neumann, J. Andrew Bagnell, Pieter Abbeel, and Jan Peters. An algorithmic perspective on imitation learning.

*Foundations and Trends® in Robotics*, 7(1-2):1–179, 2018. ISSN 1935-8253. doi: 10.1561/23000000053. URL <http://dx.doi.org/10.1561/23000000053>.

- [29] Long Ouyang, Jeff Wu, Xu Jiang, Diogo Almeida, Carroll L. Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, John Schulman, Jacob Hilton, Fraser Kelton, Luke Miller, Maddie Simens, Amanda Askell, Peter Welinder, Paul Christiano, Jan Leike, and Ryan Lowe. Training language models to follow instructions with human feedback, 2022.
- [30] Aravind Rajeswaran, Vikash Kumar, Abhishek Gupta, Giulia Vezzani, John Schulman, Emanuel Todorov, and Sergey Levine. Learning complex dexterous manipulation with deep reinforcement learning and demonstrations. In *Robotics: Science and Systems*, 2018.
- [31] Siddharth Reddy, Anca D. Dragan, and Sergey Levine. Sqil: Imitation learning via reinforcement learning with sparse rewards. *arXiv: Learning*, 2019.
- [32] Stéphane Ross and Drew Bagnell. Efficient reductions for imitation learning. In *International Conference on Artificial Intelligence and Statistics (AISTATS)*, pages 661–668, 2010.
- [33] Stéphane Ross and J. Andrew Bagnell. Reinforcement and imitation learning via interactive no-regret learning. In *Advances in Neural Information Processing Systems*, 2014.
- [34] Stéphane Ross, Geoffrey Gordon, and Drew Bagnell. A reduction of imitation learning and structured prediction to no-regret online learning. In *International Conference on Artificial Intelligence and Statistics (AISTATS)*, pages 627–635, 2011.
- [35] David Silver, Julian Schrittwieser, Karen Simonyan, Ioannis Antonoglou, Aja Huang, Arthur Guez, Thomas Hubert, Lucas Baker, Matthew Lai, Adrian Bolton, et al. Mastering the game of go without human knowledge. *nature*, 550(7676): 354–359, 2017.
- [36] Yuda Song, Yifei Zhou, Ayush Sekhari, J Andrew Bagnell, Akshay Krishnamurthy, and Wen Sun. Hybrid rl: Using both offline and online data can make rl efficient. *arXiv preprint arXiv:2210.06718*, 2022.

- [37] Jonathan Spencer, Sanjiban Choudhury, Matthew Barnes, Matthew Schmittle, Mung Chiang<sup>1</sup>, Peter Ramadge<sup>1</sup>, and Siddhartha Srinivasa<sup>2</sup>. Learning from interventions: Human-robot interaction as both explicit and implicit feedback. In *Proceedings of Robotics: Science and Systems*, 2020.
- [38] Wen Sun, Arun Venkatraman, Geoffrey J. Gordon, Byron Boots, and J. Andrew Bagnell. Deeply aggravated: Differentiable imitation learning for sequential prediction. In *Proceedings of Machine Learning Research*, 2017.
- [39] Wen Sun, J. Andrew Bagnell, and Byron Boots. Truncated horizon policy search: Combining reinforcement learning and imitation learning. 2018.
- [40] Mel Vecerik, Todd Hester, Jonathan Scholz, Fumin Wang, Olivier Pietquin, Bilal Piot, Nicolas Heess, Thomas Rothörl, Thomas Lampe, and Martin Riedmiller. Leveraging demonstrations for deep reinforcement learning on robotics problems with sparse rewards. *arXiv preprint arXiv:1707.08817*, 2017.
- [41] Tengyang Xie, Nan Jiang, Huan Wang, Caiming Xiong, and Yu Bai. Policy finetuning: Bridging sample-efficient offline and online reinforcement learning. 2022.
- [42] Zhenghai Xue, Zhenghao Peng, Quanyi Li, Zhihan Liu, and Bolei Zhou. Guarded policy optimization with imperfect online demonstrations, 2023.