

Recovery RL: Safe Reinforcement Learning with Learned Recovery Zones

Brijen Thananjeyan^{*1}, Ashwin Balakrishna^{*1}, Suraj Nair², Michael Luo¹, Krishnan Srinivasan², Minh Hwang¹, Joseph E. Gonzalez¹, Julian Ibarz³, Chelsea Finn², Ken Goldberg¹

Abstract—Safety remains a central obstacle preventing widespread use of RL in the real world: learning new tasks in uncertain environments requires extensive exploration, but safety requires limiting exploration. We propose Recovery RL, an algorithm which navigates this tradeoff by (1) leveraging offline data to learn about constraint violating zones *before* policy learning and (2) *separating* the goals of improving task performance and constraint satisfaction across two policies: a task policy that only optimizes the task reward and a recovery policy that guides the agent to safety when constraint violation is likely. We evaluate Recovery RL on 6 simulation domains, including two contact-rich manipulation tasks and an image-based navigation task, and an image-based obstacle avoidance task on a physical robot. We compare Recovery RL to 5 prior safe RL methods which jointly optimize for task performance and safety via constrained optimization or reward shaping and find that Recovery RL outperforms the next best prior method across all domains. Results suggest that Recovery RL trades off constraint violations and task successes 2 - 20 times more efficiently in simulation domains and 3 times more efficiently in physical experiments. See <https://tinyurl.com/rl-recovery> for videos and supplementary material.

Index Terms—Reinforcement Learning, Safety

I. INTRODUCTION

Reinforcement learning (RL) provides a general framework for robots to acquire new skills, and has shown promise in a variety of robotic domains such as navigation [27], locomotion [15], and manipulation [19, 23]. However, when deploying RL agents in the real world, unconstrained exploration can result in highly suboptimal behaviors which can damage the robot, break surrounding objects, or bottleneck the learning process. For example, consider an agent tasked with learning to extract a carton of milk from a fridge. If it tips over the carton, then not only can this possibly break the carton and create a mess, but it also requires laborious human effort to wipe up the milk and replace the carton so

^{*} equal contribution

Manuscript received: October, 15, 2019; Revised February 2, 2020; Accepted February 27, 2020.

This paper was recommended for publication by Editor Dana Kulic upon evaluation of the Associate Editor and Reviewers' comments.

¹Brijen Thananjeyan, Ashwin Balakrishna, Suraj Nair, Michael Luo, Minh Hwang, Joseph E. Gonzalez, and Ken Goldberg are with the Dept. of Electrical Engineering and Computer Science, University of California, Berkeley, USA {bthananjeyan, ashwin_balakrishna, michael.luo, gkgkgk1215, jegonzal, goldberg}@berkeley.edu

²Suraj Nair, Krishnan Srinivasan, and Chelsea Finn are with the Dept. of Computer Science, Stanford University, USA {surajn, krshna, cbfinn}@stanford.edu

³Julian Ibarz is with Google AI, USA julianibarz@google.com
Digital Object Identifier (DOI): see top of this page.

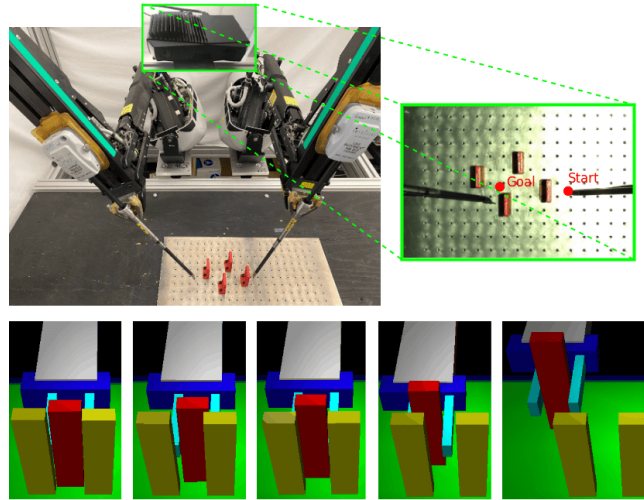


Figure 1: Recovery RL can safely learn policies for contact-rich tasks from high-dimensional image observations in simulation experiments and on a physical robotic system. We evaluate Recovery RL on an image-based obstacle avoidance task with delta-position control on the da Vinci Research Kit (top left) with overhead image observations (top right). We find that Recovery RL substantially outperforms prior methods (Figure 6), suggesting that it can be used for visuomotor control on physical robots. We also find that Recovery RL can perform challenging contact-rich manipulation tasks in simulation; as shown in the bottom row, Recovery RL successfully extracts the red block without toppling other blocks by learning to nudge it away from other blocks before grasping it.

that the robot can continue learning. In the meantime, the robot is not able to collect experience or improve its policy until the consequences of this violation are rectified. Thus, endowing RL agents with the ability to satisfy constraints during learning not only enables robots to interact safely, but also allows them to more efficiently learn in the real world. However, enforcing constraints on the agent's behavior during learning is challenging, since system dynamics and the states leading to constraint violations may be initially unknown and must be learned from experience, especially when learning from high dimensional observations such as images. Safe exploration poses a tradeoff: learning new skills through environmental interaction requires exploring a wide range of possible behaviors, but learning safely forces the agent to restrict exploration to constraint satisfying states.

We consider a RL formulation subject to constraints on the probability of unsafe future behavior and design an algorithm that can balance the often conflicting objectives of task-directed exploration and safety. Most prior work in safe RL integrates constraint satisfaction into the task objective to

jointly optimize the two. While these approaches are appealing for their generality and simplicity, there are two key aspects which make them difficult to use in practice. First, the inherent objective conflict between exploring to learn new tasks and limiting exploration to avoid constraint violations can lead to suboptimality in policy optimization. Second, exploring the environment to learn about constraints requires a significant amount of constraint violations during learning. However, this can result in the agent taking uncontrolled actions which can damage itself and the environment.

We take a step towards addressing these issues with two key algorithmic ideas. First, inspired by recent work in robust control [4, 11, 13, 21], we represent the RL agent with two policies: the first policy focuses on optimizing the unconstrained task objective (task policy) and the second policy takes control when the task policy is in danger of constraint violations in the near future (recovery policy). Instead of modifying the policy optimization procedure to encourage constraint satisfaction, which can introduce suboptimality in the learned task policy [26], the recovery policy can be viewed as defining an alternate MDP for the task policy to explore in which constraint violations are unlikely. Separating the task and recovery policies makes it easier to balance task performance and safety, and allows using off-the-shelf RL algorithms for both. Second, we leverage offline data to learn a recovery set, which indicates regions of the MDP in which future constraint violations are likely, and a recovery policy, which is queried within this set to prevent violations. This offline data can be collected under human supervision to illustrate examples of desired behaviors before the agent interacts with the environment or can contain unsafe behaviors previously experienced by the robot in the environment when performing other tasks. Both the recovery set and policy are updated online with agent experience, but the offline data allows the agent to observe constraint violations and learn from them without the task policy directly having to experience too many uncontrolled violations during learning.

We present Recovery RL, a new algorithm for safe robotic RL. Unlike prior work, Recovery RL (1) can leverage offline data of constraint violations to learn about constraints *before* interacting with the environment, and (2) uses separate policies for the task and recovery to learn safely without significantly sacrificing task performance. We evaluate Recovery RL against 5 state-of-the-art safe RL algorithms on 6 navigation and manipulation domains in simulation, including a visual navigation task, and find that Recovery RL trades off constraint violations and task successes 2 - 20 times more efficiently than the next best prior method. We evaluate Recovery RL on an image-based obstacle avoidance task on a physical robot and find that it trades off constraint violations and task successes 3 times more efficiently than the next best prior algorithm.

II. RELATED WORK

Prior work has studied safety in RL in several ways, including imposing constraints on expected return [1, 34], risk measures [17, 29, 31, 33], and avoiding regions of the MDP where constraint violations are likely [4, 5, 10, 12, 36, 37].

We build on the latter approach and design an algorithm which uses a learned recovery policy to keep the RL agent within a learned safe region of the MDP.

Jointly Optimizing for Task Performance and Safety: A popular strategy in algorithms for safe RL involves modifying the policy optimization procedure of standard RL algorithms to simultaneously reason about both task reward and constraints using methods such as trust regions [1], optimizing a Lagrangian relaxation [25, 30, 34], or constructing Lyapunov functions [6, 7]. The most similar of these works to Recovery RL is Srinivasan *et al.* [30], which trains a safety critic to estimate the probability of future constraint violation under the current policy and optimizes a Lagrangian objective function to limit the probability of constraint violations while maximizing task reward. Unlike Srinivasan *et al.* [30], which uses the safety critic to modify the task policy optimization objective, Recovery RL uses it to determine when to execute a learned recovery policy which minimizes the safety critic to keep the agent in safe regions of the MDP. This idea enables Recovery RL to more effectively balance task performance and constraint satisfaction than algorithms which jointly optimize for task performance and safety.

Restricting Exploration with an Auxiliary Policy: Another approach to safe RL explicitly restricts policy exploration to a safe subset of the MDP using a recovery or shielding mechanism. This idea has been explored in [4, 11], which utilize Hamilton-Jacobi reachability analysis to define a task policy and safety controller, and in the context of shielding [2, 13, 21]. In contrast to these works, which assume approximate knowledge of system dynamics or require precise knowledge of constraints a priori, Recovery RL learns information about the MDP, such as constraints and dynamics, from a combination of offline data and online experience. This allows Recovery RL to scale to high-dimensional state spaces such as images, in which exact specification of system dynamics and constraints can be very challenging, and is often impossible. Additionally, Recovery RL reasons about chance constraints rather than robust constraints, which may be challenging to satisfy when dynamics are unknown. Fisac *et al.* [11] design and prove safety guarantees for learning-based controllers in a robust optimal control setting with known dynamics and a robust control invariant safe set. With these additional assumptions, Recovery RL has similar theoretical properties as well. Han *et al.* [16] and Eysenbach *et al.* [10] introduce reset policies which are trained jointly with the task policy to reset the agent to its initial state distribution, ensuring that the task policy only learns behaviors which can be reset [10]. However, enforcing the ability to fully reset can be impractical or inefficient. Inspired by this work, Recovery RL instead executes approximate resets to nearby safe states when constraint violation is probable. Richter *et al.* [27] learns the probability of constraint violation conditioned on an action plan to activate a hand-designed safety controller. In contrast, Recovery RL uses a learned recovery mechanism which can be broadly applied across different tasks.

Leveraging Demonstrations for Safe RL and Control: There has also been significant prior work investigating how demonstrations can be leveraged to enable safe exploration.

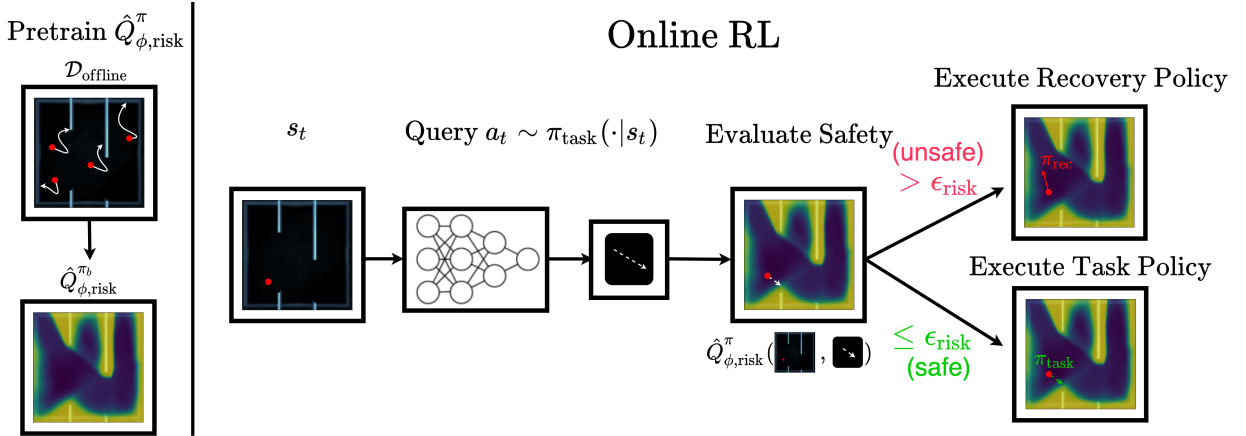


Figure 2: **Recovery RL**: For intuition, we illustrate Recovery RL on a 2D maze navigation task where a constraint violation corresponds to hitting a wall. Recovery RL first learns safety critic $\hat{Q}_{\phi, \text{risk}}^{\pi}$ with offline data from some behavioral policy π_b , which provides a small number of controlled demonstrations of constraint violating behavior as shown on the left. For the purposes of illustration, we visualize the average of the $\hat{Q}_{\phi, \text{risk}}^{\pi}$ learned by Recovery RL over 100 action samples. Then, at each timestep, Recovery RL queries the task policy π_{task} for some action a at state s , evaluates $\hat{Q}_{\phi, \text{risk}}^{\pi}(s, a)$, and executes the recovery policy π_{rec} if $\hat{Q}_{\phi, \text{risk}}^{\pi}(s, a) > \epsilon_{\text{risk}}$ and π_{task} otherwise. The task policy, recovery policy, and safety critic are updated after each transition from agent experience.

Rosolia *et al.* [28] and Thananjeyan *et al.* [35] introduce model predictive control algorithms which leverage initial constraint satisfying demonstrations to iteratively improve their performance with safety guarantees and Thananjeyan *et al.* [36] extends these ideas to the RL setting. In contrast to these works, Recovery RL learns a larger safe set that explicitly models future constraint satisfaction and also learns the problem constraints from prior experience without task specific demonstrations. Also, Recovery RL is compatible with model-free RL algorithms while [35, 36] require a dynamics model to evaluate reachability-based safety online.

III. PROBLEM STATEMENT

We consider RL under Markov decision processes (MDPs), which can be described by tuple $\mathcal{M} = (\mathcal{S}, \mathcal{A}, P(\cdot|\cdot, \cdot), R(\cdot, \cdot), \gamma, \mu)$ where \mathcal{S} and \mathcal{A} are the state and action spaces. Stochastic dynamics model $P: \mathcal{S} \times \mathcal{A} \times \mathcal{S} \rightarrow [0, 1]$ maps a state and action to a probability distribution over subsequent states, $\gamma \in [0, 1]$ is a discount factor, μ is the initial state distribution ($s_0 \sim \mu$), and $R: \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}$ is the reward function. We augment the MDP with an extra constraint cost function $C: \mathcal{S} \rightarrow \{0, 1\}$ which indicates whether a state is constraint violating and associated discount factor $\gamma_{\text{risk}} \in [0, 1]$. This yields the following new MDP: $(\mathcal{S}, \mathcal{A}, P(\cdot|\cdot, \cdot), R(\cdot, \cdot), \gamma, C(\cdot), \gamma_{\text{risk}})$. We assume that episodes terminate on violations, equivalent to transitioning to a constraint-satisfying absorbing state with zero reward.

Let Π be the set of Markovian stationary policies. Given policy $\pi \in \Pi$, the expected return is defined as $R^{\pi} = \mathbb{E}_{\pi, \mu, P} [\sum_t \gamma^t R(s_t, a_t)]$ and the expected discounted probability of constraint violation is defined as $Q_{\text{risk}}^{\pi}(s_i, a_i) = \mathbb{E}_{\pi, \mu, P} [\sum_t \gamma_{\text{risk}}^t C(s_{t+i})] = \sum_t \gamma_{\text{risk}}^t \mathbb{P}(C(s_{t+i}) = 1)$, which we would like to be below a threshold $\epsilon_{\text{risk}} \in [0, 1]$. The goal is to solve the following constrained optimization problem:

$$\pi^* = \arg \max_{\pi \in \Pi} \{R^{\pi} : Q_{\text{risk}}^{\pi}(s_0, a_0) \leq \epsilon_{\text{risk}}\} \quad (\text{III.1})$$

This setting exactly corresponds to the CMDP formulation from [3], but with constraint costs limited to binary indicator functions for constraint violating states. We limit the choice to binary indicator functions, as they are easier to provide than shaped costs and use Q_{risk}^{π} to convey information about delayed constraint costs. We define the set of feasible policies, $\{\pi : Q_{\text{risk}}^{\pi} \leq \epsilon\}$, the set of ϵ -safe policies Π_{ϵ} . Observe that if $\gamma_{\text{risk}} = 1$, then by the assumption of termination on constraint violation, $Q_{\text{risk}}^{\pi}(s_i, a_i) = \mathbb{P}(\bigcup_t C(s_t) = 1)$, or the probability of a constraint violation in the future. Setting $\epsilon_{\text{risk}} = 0$ as well results in a robust optimal control problem.

We present an algorithm to optimize equation (III.1) by utilizing a pair of policies, a *task policy* π_{task} , which is trained to maximize R^{π} over $\pi_{\text{task}} \in \Pi$ and a *recovery policy* π_{rec} , which attempts to guide the agent back to a state-action tuple (s, a) where $Q_{\text{risk}}^{\pi}(s, a) \leq \epsilon_{\text{risk}}$. We assume access to a set of transitions from offline data ($\mathcal{D}_{\text{offline}}$) with examples of constraint violations. Unlike in typical imitation learning settings, this data need not illustrate task successes, but shows possible ways to violate constraints. We leverage $\mathcal{D}_{\text{offline}}$ to constrain exploration of the task policy to reduce the probability of constraint violation during environment interaction.

IV. RECOVERY RL

We outline the central ideas behind Recovery RL. In Section IV-A, we review how to learn a safety critic to estimate the probability of future constraint violations for the agent’s policy. Then in Section IV-B, we show how this safety critic is used to define the recovery policy for Recovery RL and the recovery set in which it is activated. In Section IV-C we discuss how the safety critic and recovery policy are initialized from offline data and in Section IV-D we discuss implementation details. See Algorithm 1 and Figure 2 for further illustration of Recovery RL.

A. Preliminaries: Training a Safety Critic

As in Srinivasan *et al.* [30], Recovery RL learns a critic function Q_{risk}^{π} that estimates the discounted future probability

of constraint violation of the current policy π :

$$\begin{aligned} Q_{\text{risk}}^{\pi}(s_t, a_t) &= \mathbb{E}_{\pi} \left[\sum_{t'=t}^{\infty} \gamma_{\text{risk}}^{t'-t} c_{t'} | s_t, a_t \right] \\ &= c_t + (1 - c_t) \gamma_{\text{risk}} \mathbb{E}_{\pi} [Q_{\text{risk}}^{\pi}(s_{t+1}, a_{t+1}) | s_t, a_t]. \end{aligned} \quad (\text{IV.1})$$

Here $c_t = 1$ indicates that state s_t is constraint violating with $c_t = 0$ otherwise. Note we do not assume access to the true constraint cost function C . This is different from the standard Bellman equations to the assumption that episodes terminate when $c_t = 1$. In practice, we train a sample-based approximation $\hat{Q}_{\phi, \text{risk}}^{\pi}$, parameterized by ϕ , by approximating these equations using sampled transitions (s_t, a_t, s_{t+1}, c_t) .

We train $\hat{Q}_{\phi, \text{risk}}^{\pi}$ by minimizing the following MSE loss with respect to the target (RHS of equation IV.1).

$$\begin{aligned} J_{\text{risk}}(s_t, a_t, s_{t+1}; \phi) &= \frac{1}{2} \left(\hat{Q}_{\phi, \text{risk}}^{\pi}(s_t, a_t) - (c_t \right. \\ &\left. + (1 - c_t) \gamma_{\text{risk}} \mathbb{E}_{a_{t+1} \sim \pi(\cdot | s_{t+1})} [\hat{Q}_{\phi, \text{risk}}^{\pi}(s_{t+1}, a_{t+1})]) \right)^2 \end{aligned} \quad (\text{IV.2})$$

We use a target network to create the target values [14, 30].

B. Defining a Recovery Set and Policy

Recovery RL executes a composite policy π in the environment, which selects between a task-driven policy π_{task} and a recovery policy π_{rec} at each timestep based on whether the agent is in danger of constraint violations in the near future. To quantify this risk, we use Q_{risk}^{π} to construct a recovery set that contains state-action tuples from which π may not be able to avoid constraint violations. Then if the agent finds itself in the recovery set, it executes a learned recovery policy instead of π_{task} to navigate back to regions of the MDP that are known to be sufficiently safe. Specifically, define two complimentary sets: the safe set $\mathcal{T}_{\text{safe}}^{\pi}$ and recovery set $\mathcal{T}_{\text{rec}}^{\pi}$:

$$\begin{aligned} \mathcal{T}_{\text{safe}}^{\pi} &= \{(s, a) \in \mathcal{S} \times \mathcal{A} : Q_{\text{risk}}^{\pi}(s, a) \leq \epsilon_{\text{risk}}\} \\ \mathcal{T}_{\text{rec}}^{\pi} &= \mathcal{S} \times \mathcal{A} \setminus \mathcal{T}_{\text{safe}}^{\pi} \end{aligned}$$

We consider state-action tuple (s, a) to be safe if in state s after taking action a , executing π has a discounted probability of constraint violation less than ϵ_{risk} .

If the task policy π_{task} proposes an action $a^{\pi_{\text{task}}}$ at state s such that $(s, a^{\pi_{\text{task}}}) \notin \mathcal{T}_{\text{safe}}^{\pi}$, then a recovery action sampled from π_{rec} is executed instead of $a^{\pi_{\text{task}}}$. Thus, the recovery policy in Recovery RL can be thought of as projecting π_{task} into a safe region of the policy space in which constraint violations are unlikely. The recovery policy π_{rec} is also an RL agent, but is trained to minimize $\hat{Q}_{\phi, \text{risk}}^{\pi}(s, a)$ to reduce the risk of constraint violations under π . Let $a_t^{\pi_{\text{task}}} \sim \pi_{\text{task}}(\cdot | s_t)$ and $a_t^{\pi_{\text{rec}}} \sim \pi_{\text{rec}}(\cdot | s_t)$. Then π selects actions as follows:

$$a_t = \begin{cases} a_t^{\pi_{\text{task}}} & (s_t, a_t^{\pi_{\text{task}}}) \in \mathcal{T}_{\text{safe}}^{\pi} \\ a_t^{\pi_{\text{rec}}} & (s_t, a_t^{\pi_{\text{task}}}) \in \mathcal{T}_{\text{rec}}^{\pi} \end{cases} \quad (\text{IV.3})$$

Recovery RL filters proposed actions that are likely to lead to unsafe states, equivalent to modifying the environment that π_{task} operates in with new dynamics:

$$P_{\epsilon_{\text{risk}}}^{\pi_{\text{rec}}}(s' | s, a) = \begin{cases} P(s' | s, a) & (s, a) \in \mathcal{T}_{\text{safe}}^{\pi} \\ P(s' | s, a^{\pi_{\text{rec}}}) & (s, a) \in \mathcal{T}_{\text{rec}}^{\pi} \end{cases} \quad (\text{IV.4})$$

We train $\hat{Q}_{\phi, \text{risk}}^{\pi}$ on samples from π since π_{task} is not executed directly in the environment, but is rather filtered through π .

It is easy to see that the proposed recovery mechanism will shield the agent from regions in which constraint violations are likely if $\hat{Q}_{\phi, \text{risk}}^{\pi}$ is correct and executing π_{rec} reduces its value. However, this poses a potential concern: while the agent may be safe, how do we ensure that π_{task} can make progress in the *new* MDP defined in equation IV.4? Suppose that π_{task} proposes an unsafe action $a_t^{\pi_{\text{task}}}$ under $\hat{Q}_{\phi, \text{risk}}^{\pi}$. Then, Recovery RL executes a recovery action $a_t^{\pi_{\text{rec}}}$ and observes transition $(s_t, a_t^{\pi_{\text{rec}}}, s_{t+1}, r_t)$ in the environment. However, if π_{task} is updated with this observed transition, it will not learn to associate its proposed action ($a_t^{\pi_{\text{task}}}$) in the new MDP with r_t and s_{t+1} . As a result, π_{task} may continue to propose the same unsafe actions without realizing it is observing the result of an action sampled from π_{rec} . To address this issue, for training π_{task} , we *relabel all actions with the action proposed by π_{task}* . Thus, instead of training π_{task} with executed transitions (s_t, a_t, s_{t+1}, r_t) , π_{task} is trained with transitions $(s_t, a_t^{\pi_{\text{task}}}, s_{t+1}, r_t)$. This ties into the interpretation of defining a safe MDP with dynamics $P_{\epsilon_{\text{risk}}}^{\pi_{\text{rec}}}(s' | s, a)$ for π_{task} to act in since all transitions for training π_{task} are relabeled as if π_{task} was executed directly.

Algorithm 1 Recovery RL

Require: $\mathcal{D}_{\text{offline}}$, task horizon H , number of episodes N

- 1: Pretrain π_{rec} and $\hat{Q}_{\phi, \text{risk}}^{\pi}$ on $\mathcal{D}_{\text{offline}}$ ▷ Section IV-C
- 2: $\mathcal{D}_{\text{task}} \leftarrow \emptyset$, $\mathcal{D}_{\text{rec}} \leftarrow \mathcal{D}_{\text{offline}}$
- 3: $s_0 \leftarrow \text{env.reset}()$
- 4: **for** $i \in \{1, \dots, N\}$ **do**
- 5: **for** $t \in \{1, \dots, H\}$ **do**
- 6: **if** $c_t = 1$ or `is_terminal`(s_t) **then**
- 7: $s_t \leftarrow \text{env.reset}()$
- 8: **end if**
- 9: $a_t^{\pi_{\text{task}}} \sim \pi_{\text{task}}(\cdot | s_t)$ ▷ Query task policy
- 10: ▷ Check if task policy will be unsafe
- 11: **if** $(s_t, a_t^{\pi_{\text{task}}}) \in \mathcal{T}_{\text{rec}}^{\pi}$ **then**
- 12: $a_t \sim \pi_{\text{rec}}(\cdot | s_t)$ ▷ Select recovery policy
- 13: **else**
- 14: $a_t = a_t^{\pi_{\text{task}}}$ ▷ Select task policy
- 15: **end if**
- 16: Execute a_t
- 17: Observe s_{t+1} , $r_t = R(s_t, a_t)$, $c_t = C(s_t)$
- 18: ▷ Relabel transition
- 19: $\mathcal{D}_{\text{task}} \leftarrow \mathcal{D}_{\text{task}} \cup \{(s_t, a_t^{\pi_{\text{task}}}, s_{t+1}, r_t)\}$
- 20: $\mathcal{D}_{\text{rec}} \leftarrow \mathcal{D}_{\text{rec}} \cup \{(s_t, a_t, s_{t+1}, c_t)\}$
- 21: Train π_{task} on $\mathcal{D}_{\text{task}}$, π_{rec} on \mathcal{D}_{rec}
- 22: Train $\hat{Q}_{\phi, \text{risk}}^{\pi}$ on \mathcal{D}_{rec} ▷ Eq. IV.2
- 23: **end for**
- 24: **end for**

C. Offline Pretraining

To convey information about constraints before interaction with the environment, we provide the agent with a set of transitions $\mathcal{D}_{\text{offline}}$ that contain constraint violations for pretraining. While this requires violating constraints in the environment, this data can be collected by human defined policies or under human supervision, and thus provide the robotic agent with

examples of constraint violations without the robot having to experience too many uncontrolled examples online. We pre-train $\hat{Q}_{\phi, \text{risk}}^{\pi}$ by minimizing Equation IV.2 over offline batches sampled from $\mathcal{D}_{\text{offline}}$. We also pretrain π_{rec} using $\mathcal{D}_{\text{offline}}$. Then, π_{task} , π_{rec} , and $\hat{Q}_{\phi, \text{risk}}^{\pi}$ are all updated online using experience from the agent’s composite policy as discussed in Section IV-B and illustrated in Algorithm 1. Any RL algorithm can be used to represent π_{task} while any off-policy RL algorithm can be used to learn π_{rec} . For some environments in which exploration is challenging, we use a separate set of task demos to initialize π_{task} to expedite learning.

D. Practical Implementation

Recovery Policy: Any off-policy RL algorithm can be used to learn π_{rec} . In this paper, we explore both model-free and model-based RL algorithms to learn π_{rec} . For model-free recovery, we perform gradient descent on the safety critic $\hat{Q}_{\phi, \text{risk}}^{\pi}(s, \pi_{\text{rec}}(s))$, as in the popular off-policy RL algorithm DDPG [22]. For model-based recovery, we perform model predictive control (MPC) over a learned dynamics model f_{θ} using the safety critic as a cost function. For lower dimensional tasks, we utilize the PETS algorithm from Chua *et al.* [8] to plan over a learned stochastic dynamics model, while for tasks with visual observations, we use a VAE based latent dynamics model. **Task Policy:** We utilize the popular maximum entropy RL algorithm SAC [14] to learn π_{task} , but note that any RL algorithm could be used. Details on the implementation of both policies is in the supplement.

V. EXPERIMENTS

In the following experiments, we aim to study whether Recovery RL can (1) more effectively trade off task performance and constraint satisfaction than prior algorithms, which jointly optimize both and (2) effectively use offline data for safe RL.

Domains: We evaluate Recovery RL on a set of 6 simulation domains (Figure 3) and an image-based obstacle avoidance task on a physical robot (Figure 6). All experiments involve policy learning under state space constraints, in which a constraint violation terminates the current episode. This makes learning especially challenging, since constraint violations directly preclude further exploration. This setting is reflective of a variety of real world environments, in which constraint violations can require halting the robot due to damage to itself or its surrounding environment.

We first consider three 2D navigation domains: Navigation 1, Navigation 2, and Maze. Here, the agent only observes its position in 2D space and experiences constraint violations if it hits obstacles, walls, or workspace boundaries. We then consider three higher dimensional tasks to evaluate whether Recovery RL can be applied to contact rich manipulation tasks (Object Extraction, Object Extraction (Dynamic Obstacle)) and vision-based continuous control (Image Maze). In the object extraction environments, the goal is to extract the red block without toppling any blocks, and in the case of Object Extraction (Dynamic Obstacle), also avoiding contact with a dynamic obstacle which moves in and out of the workspace. Image Maze is a shorter horizon version of Maze, but the

agent is only provided with image observations rather than its (x, y) position in the environment.

We then evaluate Recovery RL on an image-based obstacle avoidance task on the da Vinci Research Kit (dVRK) [20] where the robot must guide its end effector within 2 mm of a target position from two possible starting locations without touching red 3D printed obstacles in the workspace. See Figure 1 for an illustration of the experimental setup. The dVRK is cable-driven and has relatively imprecise controls, motivating closed-loop control strategies to compensate for these errors [18]. Furthermore, the dVRK system has been used in the past to evaluate safe RL algorithms [36] due to its high cost and the delicate structure of its arms, which make safe learning critical. Further environment, task, and data collection details can be found in the supplement for all simulation and physical experiments.

Offline Data Collection: To effectively initialize $\hat{Q}_{\phi, \text{risk}}^{\pi}$, $\mathcal{D}_{\text{offline}}$ should ideally contain a diverse set of trajectories which violate constraints in different ways. Since $\mathcal{D}_{\text{offline}}$ need not be task specific, data from other tasks in the environment could be used, or simple human defined policies can be used to illustrate constraint violating behaviors. We take the latter approach: for all navigation environments (Navigation 1, Navigation 2, Maze, Image Maze, and the physical experiment), offline data is collected by initializing the agent in various regions of the environment and directing the agent towards the closest obstacle. For the object extraction environments (Object Extraction, Object Extraction (Dynamic Obstacle)), demonstrations are collected by guiding the end effector towards the target red block and adding Gaussian noise to controls when it is sufficiently close to the target object to make toppling likely. Recovery RL and all comparisons which have a safety critic are given the same offline dataset $\mathcal{D}_{\text{offline}}$. See the supplementary material for details on the data collection procedure, and the number of total transitions and constraint violating states for all offline datasets.

Evaluation Metric: Since Recovery RL and prior methods trade off between safety and task progress, we report the ratio of the cumulative number of task successes and the cumulative number of constraint violations at each episode to illustrate this (higher is better). We tune all algorithms to maximize this ratio, and task success is determined by defining a goal set in the state space for each environment. To avoid issues with division by zero, we add 1 to the cumulative task successes and constraint violations when computing this ratio. This metric provides a single scalar value to quantify how efficiently different algorithms balance task completion and constraint satisfaction. We do not report reward per episode, as episodes terminate on task completion or constraint violation. Each run for simulation experiments is replicated across 10 random seeds and we report the mean and standard error. For physical experiments we run each algorithm across 3 random seeds and visualize all 3 runs. In the supplementary material, we also report additional metrics for each experiment: cumulative task successes, cumulative constraint violations, and reward learning curves. We find that Recovery RL violates constraints less often than comparisons while maintaining a similar task success rate and more efficiently optimizing the task reward.

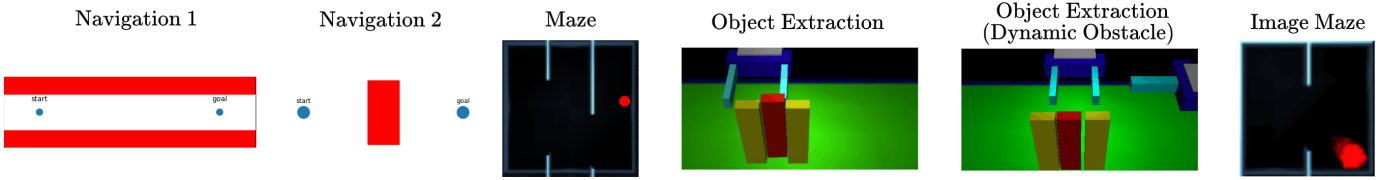


Figure 3: **Simulation Experiments Domains:** We evaluate Recovery RL on a set of 2D navigation tasks, two contact rich manipulation environments, and a visual navigation task. In Navigation 1 and 2, the goal is to navigate from the start set to the goal set without colliding into the obstacles (red) while in the Maze navigation tasks, the goal is to navigate from the left corridor to the red dot in the right corridor without colliding into walls/borders. In both object extraction environments, the objective is to grasp and lift the red block without toppling any of the blocks or colliding with the distractor arm (Dynamic Obstacle environment).

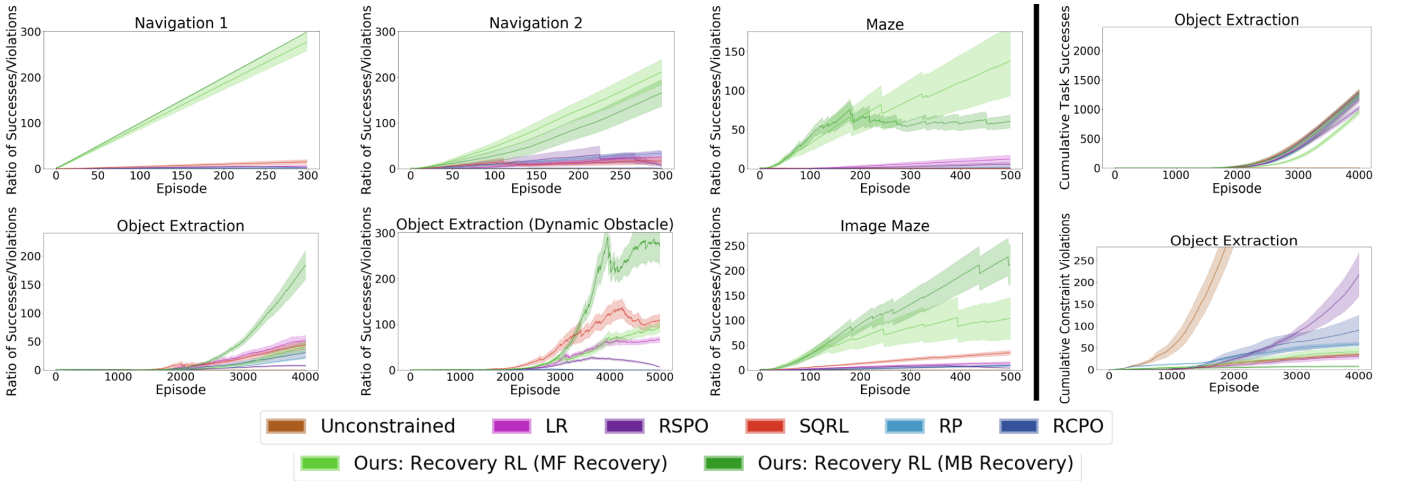


Figure 4: **Simulation Experiments:** Left: ratio of successes to constraint violations over the course of online training. In all navigation tasks, we find that Recovery RL significantly outperforms prior methods with both model-free and model-based recovery policies, while for the object extraction environments, Recovery RL with a model-based recovery policy significantly outperforms prior algorithms while Recovery RL with a model-free recovery policy does not perform as well. We hypothesize that this is due to the model-based recovery mechanism being better able to compensate for imperfections in $\hat{Q}_{\phi, \text{risk}}^{\pi}$. Results are averaged over 10 runs for each algorithm; the sawtooth pattern occurs due to constraint violations, which result in a sudden drop in the ratio. Right: cumulative successes and constraint violations. Additionally, we show the cumulative task successes and cumulative constraint violations for the Object Extraction task for all algorithms, and find that Recovery RL with model-based recovery succeeds more often than all comparisons while also violating constraints the least. Similar plots for all other experimental domains can be found in the supplementary material.

Comparisons: We compare Recovery RL to the following algorithms that ignore constraints (Unconstrained) or enforce constraints via the optimization objective (LR, SQRL, RSPO) or via reward shaping (RP, RCPO).

- **Unconstrained:** optimizes task reward, ignoring constraints.
- **Lagrangian Relaxation (LR):** minimizes $L_{\text{policy}}(s, a, r, s'; \pi) + \lambda (\mathbb{E}_{a \sim \pi(\cdot|s)} [\hat{Q}_{\phi, \text{risk}}^{\pi}(s, a)] - \epsilon_{\text{risk}})$, where L_{policy} is the policy optimization loss and the second term approximately enforces $\hat{Q}_{\phi, \text{risk}}^{\pi}(s, a) \leq \epsilon_{\text{risk}}$. Policy parameters and λ are updated via dual gradient descent.
- **Safety Q-Functions for RL (SQRL)** [30]: combines the LR method with a filtering mechanism to reject policy actions for which $\hat{Q}_{\phi, \text{risk}}^{\pi}(s, a) > \epsilon_{\text{risk}}$.
- **Risk Sensitive Policy Optimization (RSPO)** [29]: minimizes $L_{\text{policy}}(s, a, r, s'; \pi) + \lambda_t (\mathbb{E}_{a \sim \pi(\cdot|s)} [\hat{Q}_{\phi, \text{risk}}^{\pi}(s, a)] - \epsilon_{\text{risk}})$, where λ_t is a sequence which decreases to 0.
- **Reward Penalty (RP):** observes a reward function that penalizes constraint violations: $R'(s, a) = R(s, a) - \lambda C(s)$.
- **Critic Penalty Reward Constrained Policy Optimization (RCPO)** [34]: optimizes the Lagrangian relaxation via dual gradient descent and the policy

gradient trick. The policy gradient update maximizes $\mathbb{E}_{\pi} \left[\sum_{t=0}^{\infty} \gamma^t (R(s_t, a_t) - \lambda \hat{Q}_{\phi, \text{risk}}^{\pi}(s_t, a_t)) \right]$ and the multiplier update is the same as in LR.

All of these algorithms are implemented with the same base algorithm for learning the task policy (Soft Actor-Critic [14]) and all but Unconstrained and RP are modified to use the same safety critic $\hat{Q}_{\phi, \text{risk}}^{\pi}$ which is pretrained on $\mathcal{D}_{\text{offline}}$ for all methods. Thus, the key difference between Recovery RL and prior methods is how $\hat{Q}_{\phi, \text{risk}}^{\pi}$ is utilized: the comparisons use a joint objective which uses $\hat{Q}_{\phi, \text{risk}}^{\pi}$ to train a single policy that optimizes for both task performance and constraint satisfaction, while Recovery RL separates these objectives across two sub-policies. We tune all prior algorithms and report the best hyperparameter settings found on each task for the ratio-based evaluation metric. Details on Recovery RL and all comparison algorithms are in the supplement.

Results: We first study the performance of Recovery RL and prior methods in all simulation domains in Figure 4. Results suggest that Recovery RL with both model-free and model-based recovery mechanisms significantly outperform prior algorithms across all 3 2D pointmass navigation environments

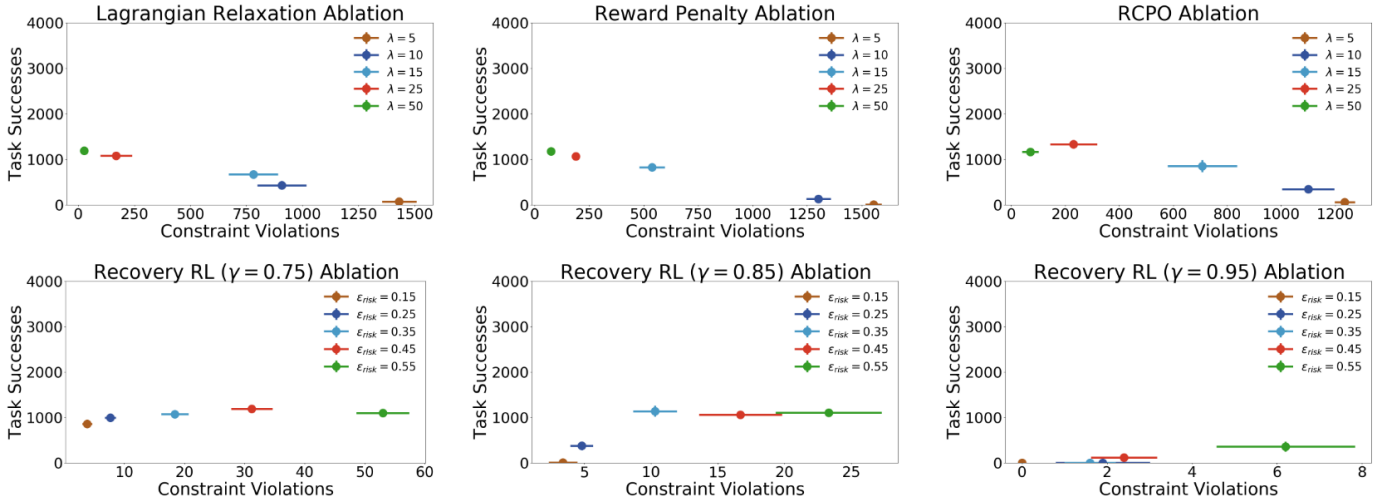


Figure 5: **Sensitivity Experiments:** We report the final number of task successes and constraint violations averaged over 10 runs at the end of training for Recovery RL and comparison algorithms for a variety of different hyperparameter settings on the Object Extraction task. We find that the comparison algorithms are relatively sensitive to the value of the penalty parameter λ while given a fixed γ_{risk} , Recovery RL achieves relatively few constraint violations while maintaining task performance over a range of ϵ_{risk} values.

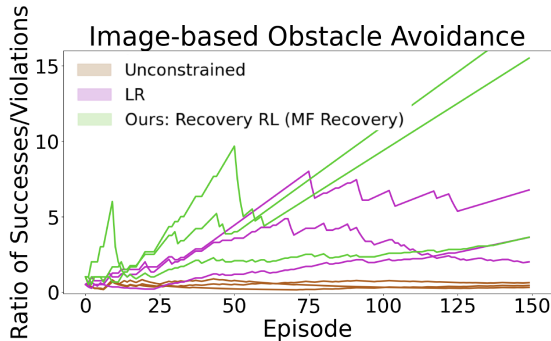


Figure 6: **Physical Experiment:** We evaluate Recovery RL on an image-based obstacle avoidance task (red obstacles) on the dVRK (Figure 1). We supply all algorithms with an overhead RGB image as input and run each algorithm 3 times. We find that Recovery RL significantly outperforms Unconstrained and LR.

(Navigation 1, Navigation 2, Maze) and the visual navigation environment (Image Maze). In the Object Extraction environments, we find that Recovery RL with model-based recovery significantly outperforms prior algorithms, while Recovery RL with a model-free recovery mechanism does not perform nearly as well. We hypothesize that the model-based recovery mechanism is better able to compensate for approximation errors in $\hat{Q}_{\phi, \text{risk}}^{\pi}$, resulting in a more robust recovery policy. We find that the prior methods often get very low ratios since they tend to achieve a similar number of task completions as Recovery RL, but with many more constraint violations. In contrast, Recovery RL is generally able to effectively trade off between task performance and safety. This is illustrated on the right pane of Figure 4, which suggests that Recovery RL with model-based recovery not only succeeds more often than comparison algorithms, but also exhibits fewer constraint violations. We study this further in the supplement. Finally, we evaluate Recovery RL and prior algorithms on the image-based obstacle avoidance task illustrated in Figure 1 and find that Recovery RL substantially outperforms prior methods, suggesting that Recovery RL can be used for contact-rich visuomotor control tasks in the real world (Figure 6). We study

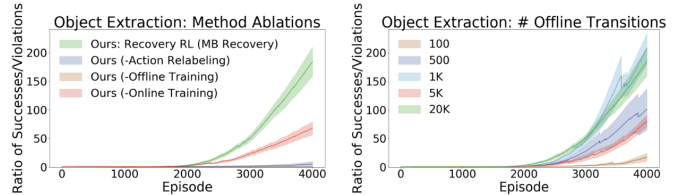


Figure 7: **Ablations:** We first study the affect of different algorithmic components of Recovery RL (left). Results suggest that offline pretraining of π_{rec} and $\hat{Q}_{\phi, \text{risk}}^{\pi}$ is critical for good performance, while removing online updates leads to a much smaller reduction in performance. Furthermore, we find that the action relabeling method for training π_{task} (Section IV-B) is critical for good performance. We then study the sensitivity of Recovery RL to the number of offline transitions used to pretrain π_{rec} and $\hat{Q}_{\phi, \text{risk}}^{\pi}$ (right) and find that Recovery RL performs well even with just 1000 transitions in $\mathcal{D}_{\text{offline}}$ for the Object Extraction task, with performance degrading when the number of transitions is reduced beyond this point.

when Recovery RL violates constraints in the supplement, and find that in most tasks, the recovery policy is already activated when constraint violations occur. This is encouraging, because if a recovery policy is challenging to learn, Recovery RL could still be used to query a human supervisor for interventions.

Ablations: We ablate different components of Recovery RL and study the sensitivity of Recovery RL to the number of transitions in $\mathcal{D}_{\text{offline}}$ for the Object Extraction domain in Figure 7. Results suggest that Recovery RL performs much more poorly when π_{rec} and $\hat{Q}_{\phi, \text{risk}}^{\pi}$ are not pretrained with data from $\mathcal{D}_{\text{offline}}$, indicating the value of learning to reason about safety before environment interaction. However, when π_{rec} and $\hat{Q}_{\phi, \text{risk}}^{\pi}$ are not updated online, performance degrades much less significantly. A key component of Recovery RL is relabeling actions when training the task policy so that π_{task} can learn to associate its proposed actions with their outcome (Section IV-B). We find that without this relabeling, Recovery RL achieves very poor performance as it rarely achieves task successes. Additionally, we find that although the reported simulation experiments supply Recovery RL and all prior methods with 20,000 transitions in $\mathcal{D}_{\text{offline}}$ for

the Object Extraction task, Recovery RL is able to achieve good performance with just 1000 transitions in $\mathcal{D}_{\text{offline}}$, with performance significantly degrading only when the size of $\mathcal{D}_{\text{offline}}$ is reduced to less than this amount.

Sensitivity Experiments: We tune hyperparameters for Recovery RL and all baselines to ensure a fair comparison. We first tune γ_{risk} and ϵ_{risk} for Recovery RL, and then use the same γ_{risk} and ϵ_{risk} for prior methods to ensure that all algorithms use the same safety critic training procedure. These two hyperparameters are the only ones tuned for Recovery RL and SQRL. For RP, RCPO, and LR, we tune the penalty term λ with γ_{risk} and ϵ_{risk} fixed as mentioned above. For RSPO, we utilize a schedule which decays λ from 2 times the best value found for λ when tuning the LR comparison to 0 with an evenly spaced linear schedule over all training episodes. In Figure 5, we study the sensitivity of Recovery RL with model-based recovery and the RP, RCPO, and LR comparisons to different hyperparameter choices on the Object Extraction task. Recovery RL appears less sensitive to hyperparameters than the comparisons for the γ_{risk} values we consider.

VI. CONCLUSION

We present Recovery RL, a new algorithm for safe RL which is able to more effectively balance task performance and constraint satisfaction than 5 state-of-the-art prior algorithms for safe RL across 6 simulation domains and an image-based obstacle avoidance task on a physical robot. In future work we hope to explore further evaluation on physical robots, establish formal guarantees, and use ideas from offline RL to more effectively pretrain the recovery policy. We will explore settings in which constraint violations may not be catastrophic and applications for large-scale robot learning.

VII. ACKNOWLEDGMENTS

This research was performed at the AUTOLAB at UC Berkeley in affiliation with the Berkeley AI Research (BAIR) Lab, the Real-Time Intelligent Secure Execution (RISE) Lab, Google Brain Robotics, and the Stanford AI Research Lab. Authors were also supported by the SAIL-Toyota Research initiative, the Scalable Collaborative Human-Robot Learning (SCHoOL) Project, the NSF National Robotics Initiative Award 1734633, and in part by donations from Google, Siemens, Amazon Robotics, Toyota Research Institute, and by equipment grants from NVIDIA. This article solely reflects the opinions and conclusions of its authors and not the views of the sponsors or their associated entities. A.B. and S.N. are supported by NSF GRFPs. We thank our colleagues who provided helpful feedback, code, and suggestions, in particular Jie Tan, Jeffrey Ichnowski, and Wisdom Agboh.

REFERENCES

- [1] J. Achiam, D. Held, A. Tamar, and P. Abbeel, "Constrained policy optimization", in *Proc. Int. Conf. on Machine Learning*, 2017.
- [2] M. Alshiekh, R. Bloem, R. Ehlers, B. Könighofer, S. Niekum, and U. Topcu, "Safe reinforcement learning via shielding", 2018.
- [3] E. Altman, *Constrained Markov Decision Processes*. 1999, p. 260.
- [4] S. Bansal, M. Chen, S. Herbert, and C. J. Tomlin, "Hamilton-jacobi reachability: A brief overview and recent advances", in *Conference on Decision and Control (CDC)*, 2017.
- [5] F. Berkenkamp, M. Turchetta, A. P. Schoellig, and A. Krause, "Safe model-based reinforcement learning with stability guarantees", in *Proc. Advances in Neural Information Processing Systems*, 2017.
- [6] Y. Chow, O. Nachum, E. Duéñez-Guzmán, and M. Ghavamzadeh, "A lyapunov-based approach to safe reinforcement learning", in *NeurIPS*, 2018.
- [7] Y. Chow, O. Nachum, A. Faust, M. Ghavamzadeh, and E. Duéñez-Guzmán, "Lyapunov-based safe policy optimization for continuous control", in *ICML Workshop RLRealLife*, 2019.
- [8] K. Chua, R. Calandra, R. McAllister, and S. Levine, "Deep reinforcement learning in a handful of trials using probabilistic dynamics models", *Proc. Advances in Neural Information Processing Systems*, 2018.
- [9] F. Ebert, C. Finn, S. Dasari, A. Xie, A. Lee, and S. Levine, "Visual foresight: Model-based deep reinforcement learning for vision-based robotic control", *arXiv preprint arXiv:1812.00568*, 2018.
- [10] B. Eysenbach, S. Gu, J. Ibarz, and S. Levine, "Leave no trace: Learning to reset for safe and autonomous reinforcement learning", *International Conference on Learning Representations*, 2018.
- [11] J. F. Fisac, A. K. Akametalu, M. N. Zeilinger, S. Kaynama, J. Gillula, and C. J. Tomlin, "A general safety framework for learning-based control in uncertain robotic systems", in *IEEE Transactions on Automatic Control*, 2018.
- [12] J. F. Fisac, N. F. Lugovoy, V. Rubies-Royo, S. Ghosh, and C. J. Tomlin, "Bridging hamilton-jacobi safety analysis and reinforcement learning", in *Proc. IEEE Int. Conf. Robotics and Automation (ICRA)*, 2019.
- [13] J. H. Gillula and C. J. Tomlin, "Guaranteed safe online learning via reachability: Tracking a ground target using a quadrotor", in *Proc. IEEE Int. Conf. Robotics and Automation (ICRA)*, 2012.
- [14] T. Haarnoja, A. Zhou, P. Abbeel, and S. Levine, "Soft actor-critic: Off-policy maximum entropy deep reinforcement learning with a stochastic actor", *Proc. Int. Conf. on Machine Learning*, 2018.
- [15] T. Haarnoja, A. Zhou, K. Hartikainen, G. Tucker, S. Ha, J. Tan, V. Kumar, H. Zhu, A. Gupta, P. Abbeel, and S. Levine, "Soft actor-critic algorithms and applications", *CoRR*, 2018.
- [16] W. Han, S. Levine, and P. Abbeel, "Learning compound multi-step controllers under unknown dynamics", *Proc. IEEE/RSJ Int. Conf. on Intelligent Robots and Systems (IROS)*, 2015.
- [17] M. Heger, "Consideration of risk in reinforcement learning", in *Machine Learning Proceedings*, 1994.
- [18] M. Hwang, B. Thananjeyan, S. Paradis, D. Seita, J. Ichnowski, D. Fer, T. Low, and K. Goldberg, "Efficiently calibrating cable-driven surgical robots with rgbd sensing, temporal windowing, and linear and recurrent neural network compensation", *Robotics and Automation Letters (RAL)*, 2020.
- [19] D. Kalashnikov, A. Irpan, P. Pastor, J. Ibarz, A. Herzog, E. Jang, D. Quillen, E. Holly, M. Kalakrishnan, V. Vanhoucke, and S. Levine, "Qt-opt: Scalable deep reinforcement learning for vision-based robotic manipulation", *Conference on Robot Learning (CoRL)*, 2018.
- [20] P. Kazanzides, Z. Chen, A. Deguet, G. S. Fischer, R. H. Taylor, and S. P. DiMaio, "An open-source research kit for the da Vinci surgical system", in *Proc. IEEE Int. Conf. Robotics and Automation (ICRA)*, 2014.
- [21] S. Li and O. Bastani, "Robust model predictive shielding for safe reinforcement learning with stochastic dynamics", in *Proc. IEEE Int. Conf. Robotics and Automation (ICRA)*, 2020.
- [22] T. P. Lillicrap, J. J. Hunt, A. Pritzel, N. Heess, T. Erez, Y. Tassa, D. Silver, and D. Wierstra, "Continuous control with deep reinforcement learning", *Proc. Int. Conf. on Learning Representations*, 2016.
- [23] A. Nagabandi, K. Konoglie, S. Levine, and V. Kumar, "Deep dynamics models for learning dexterous manipulation", *Conference on Robot Learning (CoRL)*, 2019.
- [24] S. Nair, S. Savarese, and C. Finn, "Goal-aware prediction: Learning to model what matters", *Proc. Int. Conf. on Machine Learning*, 2020.
- [25] F. W. Peterr Geibel, "Risk-sensitive reinforcement learning applied to control under constraints", in *Journal of Artificial Intelligence Research*, vol. 24, 2005.
- [26] A. Ray, J. Achiam, and D. Amodei, "Benchmarking safe exploration in deep reinforcement learning", in *NeurIPS Deep Reinforcement Learning Workshop*, 2019.
- [27] C. Richter and N. Roy, "Safe visual navigation via deep learning and novelty detection", *Robotics Science and Systems (RSS)*, 2013.
- [28] U. Rosolia and F. Borrelli, "Learning model predictive control for iterative tasks. a data-driven control framework", *IEEE Transactions on Automatic Control*, 2018.
- [29] Y. Shen, M. J. Tobia, T. Sommer, and K. Obermayer, "Risk-sensitive reinforcement learning", in *Neural Computation*, vol. 26, 2014.
- [30] K. Srinivasan, B. Eysenbach, S. Ha, J. Tan, and C. Finn, "Learning to be safe: Deep rl with a safety critic", *arXiv preprint arXiv:2010.14603*, 2020.
- [31] A. Tamar, Y. Glassner, and S. Mannor, "Policy gradients beyond expectations: Conditional value-at-risk", in *CoRR*, 2014.
- [32] P. Tandon, *Pytorch implementation of soft actor critic*, <https://github.com/pranz24/pytorch-soft-actor-critic>, 2020.
- [33] Y. C. Tang, J. Zhang, and R. Salakhutdinov, "Worst cases policy gradients", *Conf. on Robot Learning (CoRL)*, 2019.
- [34] C. Tessler, D. J. Mankowitz, and S. Mannor, "Reward constrained policy optimization", in *Proc. Int. Conf. on Learning Representations*, 2019.
- [35] B. Thananjeyan, A. Balakrishna, U. Rosolia, J. E. Gonzalez, A. Ames, and K. Goldberg, "Abc-mpc: Safe sample-based learning mpc for stochastic nonlinear dynamical systems with adjustable boundary conditions", in *Workshop on the Algorithmic Foundations of Robotics*, 2020.
- [36] B. Thananjeyan, A. Balakrishna, U. Rosolia, F. Li, R. McAllister, J. E. Gonzalez, S. Levine, F. Borrelli, and K. Goldberg, "Safety augmented value estimation from demonstrations (saved): Safe deep

- model-based rl for sparse cost robotic tasks”, *Robotics and Automation Letters (RAL)*, 2020.
- [37] M. Turchetta, F. Berkenkamp, and A. Krause, “Safe exploration in finite markov decision processes with gaussian processes”, in *Proc. Advances in Neural Information Processing Systems*, 2016.
- [38] Q. Vuong, *Pytorch implementation of pets*, <https://github.com/quanvuong/handful-of-trials-pytorch>, 2020.

Recovery RL: Safe Reinforcement Learning with Learned Recovery Zones Supplementary Material

The supplementary material is structured as follows: In Section VIII we discuss brief theoretical motivation for Recovery RL and possible variants and in Section IX we discuss algorithmic details for Recovery RL and comparison algorithms. In Section X, we report additional metrics for all domains and comparisons and in Section XI we visualize the safety critic for all navigation experiments. We provide additional details about algorithm implementation in Section XII, and on domain-specific algorithm hyperparameters in Section XIII. Finally, we report simulation and physical environment details in Section XIV.

VIII. RECOVERY RL THEORETICAL MOTIVATION AND VARIANTS

In this section, we will briefly and informally discuss additional properties of Recovery RL and then discuss some variants of Recovery RL.

A. Theoretical Motivation

Recall that the task policy is operating in an environment with modified dynamics:

$$P_{\text{risk}}^{\pi_{\text{rec}}}(s'|s, a) = \begin{cases} P(s'|s, a) & (s, a) \in \mathcal{T}_{\text{safe}}^{\pi} \\ P(s'|s, a^{\pi_{\text{rec}}}) & (s, a) \in \mathcal{T}_{\text{rec}}^{\pi} \end{cases} \quad (\text{VIII.1})$$

However, $P_{\text{risk}}^{\pi_{\text{rec}}}$ changes over time (even within the same episode) and analysis of policy learning in non-stationary MDPs is currently challenging and ongoing work. Assuming that $P_{\text{risk}}^{\pi_{\text{rec}}}$ is stationary following the pretraining phase, it is immediate that π_{task} is operating in a stationary MDP $\mathcal{M}' = (\mathcal{S}, \mathcal{A}, P_{\text{risk}}^{\pi_{\text{rec}}}, R(\cdot, \cdot), \gamma, \mu)$, and therefore all properties of π_{task} in stationary MDPs apply in \mathcal{M}' . Observe that iterative improvement for π_{task} in \mathcal{M}' implies iterative improvement for π in \mathcal{M} , since both MDPs share the same reward function, and an action taken by π_{task} in \mathcal{M}' is equivalent to π_{task} trying the action in \mathcal{M} before being potentially caught by π_{rec} .

B. Safety Value Function

One variant of Recovery RL can use a safety critic that is a state-value function $V_{\text{risk}}^{\pi}(s)$ instead of a state-action-value function. While this implementation is simpler, the Q_{risk}^{π} version used in the paper can switch to a safe action instead of an unsafe one instead of waiting to reach an unsafe state to start recovery behavior.

C. Reachability-based Variant

Another variant can use the learned dynamics model in the model-based recovery policy to perform a one (or k) step lookahead to see if future states-action tuples are in $\mathcal{T}_{\text{safe}}^{\pi}$. While Q_{risk}^{π} in principle carries information about future safety, this is an alternative method to check future states.

IX. ALGORITHM DETAILS

A. Recovery RL

Recovery Policy: In principle, any off-policy reinforcement learning algorithm can be used to learn the recovery policy π_{rec} . In this paper, we explore both model-free and model-based reinforcement learning algorithms to learn π_{rec} . For model-free recovery, we perform gradient descent on the safety critic $\hat{Q}_{\phi, \text{risk}}^{\pi}(s, \pi_{\text{rec}}(s))$, as in the popular off-policy reinforcement learning algorithm DDPG [22]. We choose the DDPG-style objective function over alternatives since we do not wish the recovery policy to explore widely. For model-based recovery, we perform model predictive control (MPC) over a learned dynamics model f_{θ} by minimizing the following objective:

$$L_{\theta}(s_r, a_r) = \mathbb{E} \left[\sum_{i=0}^H \hat{Q}_{\phi, \text{risk}}^{\pi}(\hat{s}_{r+i}, a_{r+i}) \right] \quad (\text{IX.1})$$

where $\hat{s}_{r+i+1} \sim f_{\theta}(\hat{s}_{r+i}, a_{r+i})$, $\hat{s}_r = s_r$, and $\hat{a}_r = a_r$. For lower dimensional tasks, we utilize the PETS algorithm from Chua *et al.* [8] to plan over a learned stochastic dynamics model while for tasks with visual observations, we utilize a VAE based latent dynamics model. In the offline pretraining phase, when model-free recovery is used, batches are sampled sequentially from $\mathcal{D}_{\text{offline}}$ and each batch is used to (1) train $\hat{Q}_{\phi, \text{risk}}^{\pi}$ and (2) optimize the DDPG policy to minimize the current $\hat{Q}_{\phi, \text{risk}}^{\pi}$. When model-based recovery is used, the data

in $\mathcal{D}_{\text{offline}}$ is first used to learn dynamics model f_{θ} using either PETS (low dimensional tasks) or latent space dynamics (image-based tasks). Then, $\hat{Q}_{\phi, \text{risk}}^{\pi}$ is separately optimized to over batches sampled from $\mathcal{D}_{\text{offline}}$. During the online RL phase, all methods are updated online using on-policy data from composite policy π .

Task Policy: In experiments, we utilize the popular maximum entropy RL algorithm SAC [14] to learn π_{task} , but note that any RL algorithm could be used to train π_{task} . In general π_{task} is only updated in the online RL phase. However, in certain domains where exploration is challenging, we pre-train SAC on a small set of task-specific demonstrations to expedite learning. To do this, like for training the model-free recovery policy, we sample batches sequentially from $\mathcal{D}_{\text{offline}}$ and each batch is used to (1) train $\hat{Q}_{\phi, \text{risk}}^{\pi}$ and (2) optimize the SAC policy to minimize the current $\hat{Q}_{\phi, \text{risk}}^{\pi}$. To ensure that π_{task} learns which actions result in recovery behavior, we train π_{task} on transitions $(s_t, a_t^{\pi_{\text{task}}}, s_{t+1})$ even if π_{rec} was executed.

B. Unconstrained

We use an implementation of the popular model-free reinforcement learning algorithm Soft Actor Critic [14, 32], which maximizes a combination of task reward and policy entropy with a stochastic actor function.

C. Lagrangian Relaxation (LR)

In this section we will briefly motivate and derive the Lagrangian relaxation comparison. As before, we desire to solve the following constrained optimization problem:

$$\min_{\pi} L_{\text{policy}}(s; \pi) \text{ s.t. } \mathbb{E}_{a \sim \pi(\cdot|s)} [Q_{\text{risk}}^{\pi}(s, a)] \leq \epsilon_{\text{risk}}$$

where L_{policy} is a policy loss function we would like to minimize (e.g. from SAC). As in prior work in solving constrained optimization problems, we can solve the following unconstrained problem instead:

$$\max_{\lambda \geq 0} \min_{\pi} L_{\text{policy}}(s; \pi) + \lambda (\mathbb{E}_{a \sim \pi(\cdot|s)} [Q_{\text{risk}}^{\pi}(s, a)] - \epsilon_{\text{risk}})$$

We aim to find a saddle point of the Lagrangian function via dual gradient descent. In practice, we use samples to approximate the expectation in the objective by sampling an action from $\pi(\cdot|s)$ each time the objective function is evaluated.

D. Risk Sensitive Policy Optimization (RSPO)

We implement Risk Sensitive Policy Optimization by implementing the Lagrangian Relaxation method as discussed in Section IX-C with a sequence of multipliers which decrease over time. This encourages initial constraint satisfaction followed by gradual increase in prioritization of the task objective and is inspired by the Risk Sensitive Q-learning algorithm from [29].

E. Safety Q-Functions for Reinforcement Learning (SQRL)

This comparison is identical to LR, except it additionally adds a Q-filter, that performs rejection sampling on the policy's distribution $\pi(\cdot|s_t)$ until it finds an action a_t such that $Q_{\text{risk}}^{\pi}(s_t, a_t) \leq \epsilon_{\text{risk}}$.

F. Reward Penalty (RP)

The reward penalty comparison simply involves subtracting a constant penalty λ from the task reward function when a constraint is violated. This is the only comparison algorithm other than Unconstrained which does not use the learned Q_{risk}^{π} or the constraint demos, but is included due to its surprising efficacy and simplicity.

G. Off Policy Reward Constrained Policy Optimization (RCPO)

In on-policy RCPO [34], the policy is optimized via policy gradient estimators by maximizing $\mathbb{E}_{\pi} [\sum_{t=0}^{\infty} (\gamma^t R(s, a) - \lambda \gamma^t \mathcal{Y}_{\text{risk}} D(s, a))]$. In this work, we use $D(s, a) = Q_{\text{risk}}^{\pi}(s, a)$ and update the Lagrange multiplier λ as in LR. We could also use $D(s, a) = C(s)$, which would be almost identical to the RP comparison. Instead of optimizing this with on-policy RL, we use SAC to optimize it in an off-policy fashion to be consistent with the other comparisons.

X. ADDITIONAL EXPERIMENTAL METRICS

In Figure 8 and Figure 9, we report cumulative task successes and constraint violations for all methods for all simulation experiments. We report these statistics for the image-based obstacle avoidance physical experiment in Figure 11. We observe that Recovery RL is generally very successful across most domains with relatively few violations. Some more successful comparisons tend to have many more constraint violations.

Additionally, in Figures 10 and 12, we plot the cumulative reward attained by the agent for Recovery RL and all comparison algorithms to evaluate whether Recovery RL learns more efficiently than comparisons while also learning safely. For all plots, we show total reward attained in each episode smoothed over a 100 episode length window. Additionally, we do not show results when constraints are violated to prevent negative bias for algorithms which may violate constraints very frequently, which accounts for gaps in the plot, especially for the unconstrained algorithm which tends to violate constraints very frequently. Thus, the plots illustrate the attained reward for all algorithms conditioned on not violating constraints, which provides a good measure on the quality of solutions found. We find that in addition to the safe learning shown by Recovery RL as evidenced by the results in Figure 9, the results in Figure 10 and Figure 12 indicate that when Recovery RL satisfies constraints, it generally converges to higher quality solutions more quickly compared to the comparison algorithms. These results provide further evidence that the way Recovery RL separates the often conflicting objectives of task directed optimization and constraint satisfaction allows it to not only be safer during learning, but also learn more efficiently.

In Table I, we report empirical results for when constraint violations occur in Table I. Results suggest that in most tasks, the recovery policy is already activated when violations do occur. Thus, in these failure cases, Recovery RL is able to successfully predict future violations, but is not able to prevent them. This is encouraging, and suggests that for environments in which a recovery policy is very challenging to learn, Recovery RL could still be used to query a human supervisor for interventions.

XI. SAFETY CRITIC VISUALIZATIONS

We visualize the safety critic after pretraining for the navigation domains in Figure 14 and observe that increasing γ_{risk} results in a more gradual increase in regions near obstacles. Increasing γ_{risk} carries more information about possible future violations in $Q_{\text{risk}}^{\pi}(s, a)$. However, increasing γ_{risk} too much causes the safety critic to bleed too much throughout the state-action space as in the right-most column, making it difficult to distinguish between safe and unsafe states.

XII. IMPLEMENTATION DETAILS

Here we overview implementation and hyperparameter details for Recovery RL and all comparisons. The recovery policy (π_{rec}) and task policy (π_{task}) are instantiated and trained in both the offline phase, in which data from $\mathcal{D}_{\text{offline}}$ is used to pre-train the recovery policy, and the online phase, in which Recovery RL updates the task policy with its exploration constrained by the learned safety critic and recovery policy. The safety critic and recovery policy are also updated online.

For all experiments, we build on the PyTorch implementation of Soft Actor Critic [15] provided in [32] and all trained networks are optimized with the Adam optimizer with a learning rate of $3e-4$. We first overview the hyperparameters and training details shared across Recovery RL and comparisons in Section XII-B and then discuss the implementation of the recovery policy for Recovery RL in Section XII-C.

A. Network Architectures

For low dimensional experiments, we represent the critic with a fully connected neural network with 2 hidden layers of size 256 each with ReLU activations. The policy is also represented with a fully connected network with 2 hidden layers of size 256 each, uses ReLU activations, and outputs the parameters of a conditional Gaussian. We use a deterministic version of the same policy for the model-free recovery policy. For image-based experiments, we represent the critic with a convolutional neural network with 3 convolutional layers to embed the input image and 2 fully connected layers to embed the input action. Then, these embeddings are concatenated and fed through two more fully connected layers. All fully connected layers have 256 hidden units each. We utilize 3 convolutional layers, with 128, 64, and 16 filters respectively. All layers utilize a kernel size of 3, stride of 2, and padding of 1. ReLU activations are used between all layers, and batch normalization units are added for the convolutional layers. For all algorithms which utilize a safety critic (Recovery RL, LR, SQRL, RSPO, RCPO), Q_{risk}^{π} is represented with the same architecture as the task critic except that a sigmoid

activation is added at the output head to ensure that outputs are on $[0, 1]$ in order to effectively learn the probability of constraint violation. The task and model-free recovery policies also use the same architectures for image-based experiments, except that they output the parameters of a conditional Gaussian over the action space or an action, respectively.

B. Global Training Details

To prevent overestimation bias, we train two copies of all critic networks to compute a pessimistic (min for task critic, max for safety critic) estimate of the Q-values. Each critic is associated with a target network, and Polyak averaging is used to smoothly anneal the parameters of the target network. We use a replay buffer of size 1000000 and target smoothing coefficient $\tau = 0.005$ for all experiments except for the manipulation environments, in which a replay buffer of size 100000 and target smoothing coefficient $\tau = 0.0002$. All networks are trained with batches of 256 transitions. Finally, for SAC we utilize entropy regularization coefficient $\alpha = 0.2$ and do not update it online. We take a gradient step with batch size 1000 to update the safety critic after each timestep. We also update the model free recovery policy if applicable with the same batch at each timestep. If using a model-based recovery policy, we update it for 5 epochs at the end of each episode. For pretraining, we train the safety critic and model-free recovery policy for 10,000 steps. We train the model-based recovery policy for 50 epochs.

C. Recovery Policy Training Details

In this section, we describe the neural network architectures and training procedures used by the recovery policies for all tasks.

1) *Model-Free Recovery*: The model-free recovery policy uses the same architecture as the task policy for all tasks, as described in Section XII-A. However, it directly outputs an action in the action space instead of a distribution over the action space and greedily minimizes $\hat{Q}_{\phi, \text{risk}}^{\pi}$ rather than including an entropy regularization term as in [14]. The recovery policy is trained at each timestep on a batch of 1000 samples from the replay buffer.

2) *Model-Based Recovery Training Details*: For the non-image-based model-based recovery policy, we use PETS [8, 38], which trains and plans over a probabilistic ensemble of neural networks. We use an ensemble of 5 neural networks with 3 hidden layers of size 200 and swish activations (except at the output layer) to output the parameters of a conditional Gaussian distribution. We use the TS- ∞ trajectory sampling scheme from Chua *et al.* [8] and optimize the MPC optimization problem with 400 samples, 40 elites, and 5 iterations for all environments. For image-based tasks, we utilize a VAE based latent dynamics model as in Nair *et al.* [24]. We train the encoder, decoder, and dynamics model jointly where the encoder and decoder and convolutional neural networks and the forward dynamics model is a fully connected network. We follow the same architecture as in Nair *et al.* [24]. For the encoder we utilize the following convolutional layers (channels, kernel size, stride): [(32, 4, 2), (32, 3, 1), (64, 4, 2), (64, 3, 1), (128, 4, 2), (128, 3, 1), (256, 4, 2), (256, 3, 1)] followed by fully connected layers of size [1024, 512, 2L] where L is the size of the latent space (predict mean and variance). All layers use ReLU activations except for the last layer. The decoder takes a sample from the latent space of dimension L and then feeds this through fully connected layers [128, 128, 128] which is followed by de-convolutional layers (channels, kernel size, stride): [(128, 5, 2), (64, 5, 2), (32, 6, 2), (3, 6, 2)]. All layers again use ReLU activations except for the last layer, which uses a Sigmoid activation. For the forward dynamics model, we use a fully connected network with layers [128, 128, 128, L] with ReLU activations on all but the final layer.

XIII. ENVIRONMENT SPECIFIC ALGORITHM PARAMETERS

We use the same γ_{risk} and ϵ_{risk} for LR, RSPO, SQRL, and RCPO. For LR, RSPO, and SQRL, we find that the initial choice of λ strongly affects the overall performance of this algorithm and heavily tune this. We use the same values of λ for LR and SQRL, and use twice the best value found for LR in as an initialization for the λ -schedule in RSPO. We also heavily tune λ for RP and RCPO. These values are shown for each environment in Tables III and II.

XIV. ENVIRONMENT DETAILS

In this section, we provide additional details about each of the environments used for evaluation.

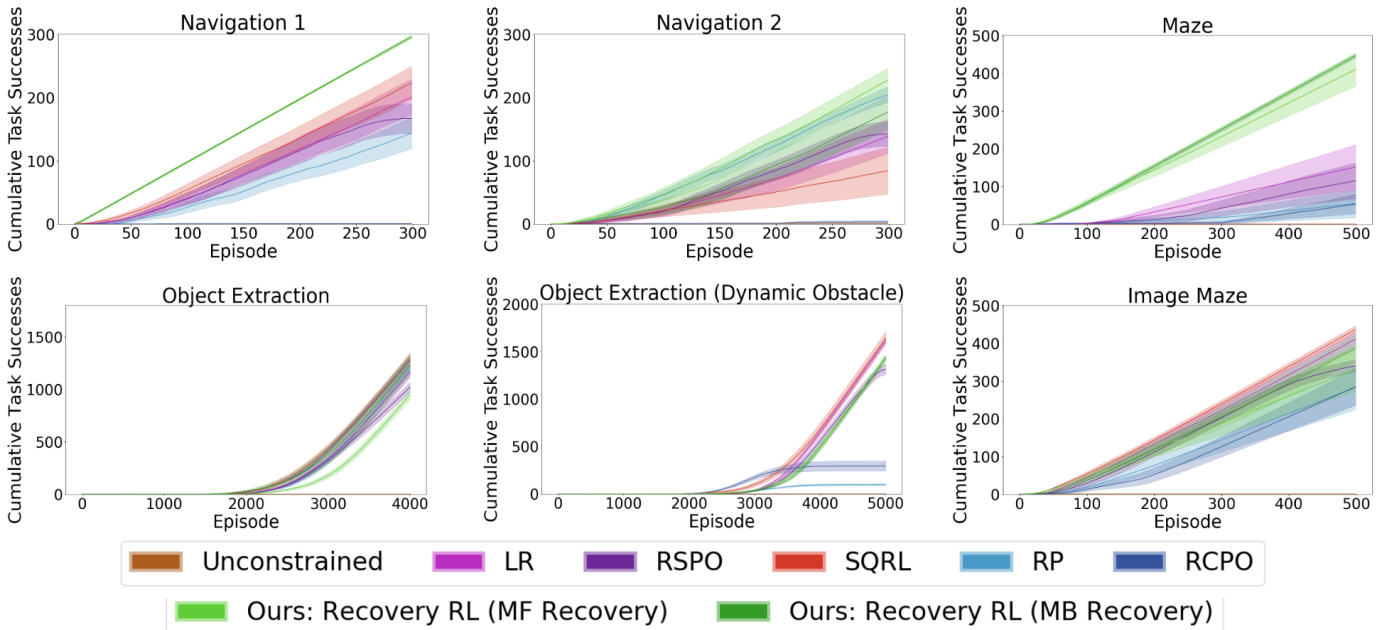


Figure 8: **Simulation Experiments Cumulative Successes:** We plot the cumulative task successes for each algorithm in each simulation domain, with results averaged over 10 runs for all algorithms. We observe that Recovery RL (green), is generally among the most successful algorithms. In the cases that it has lower successes, we observe that it is safer (Figure 9). We find that Recovery RL has a higher or comparable task success rate to the next best algorithm on all environments except for the Object Extraction (Dynamic Obstacle) environment.

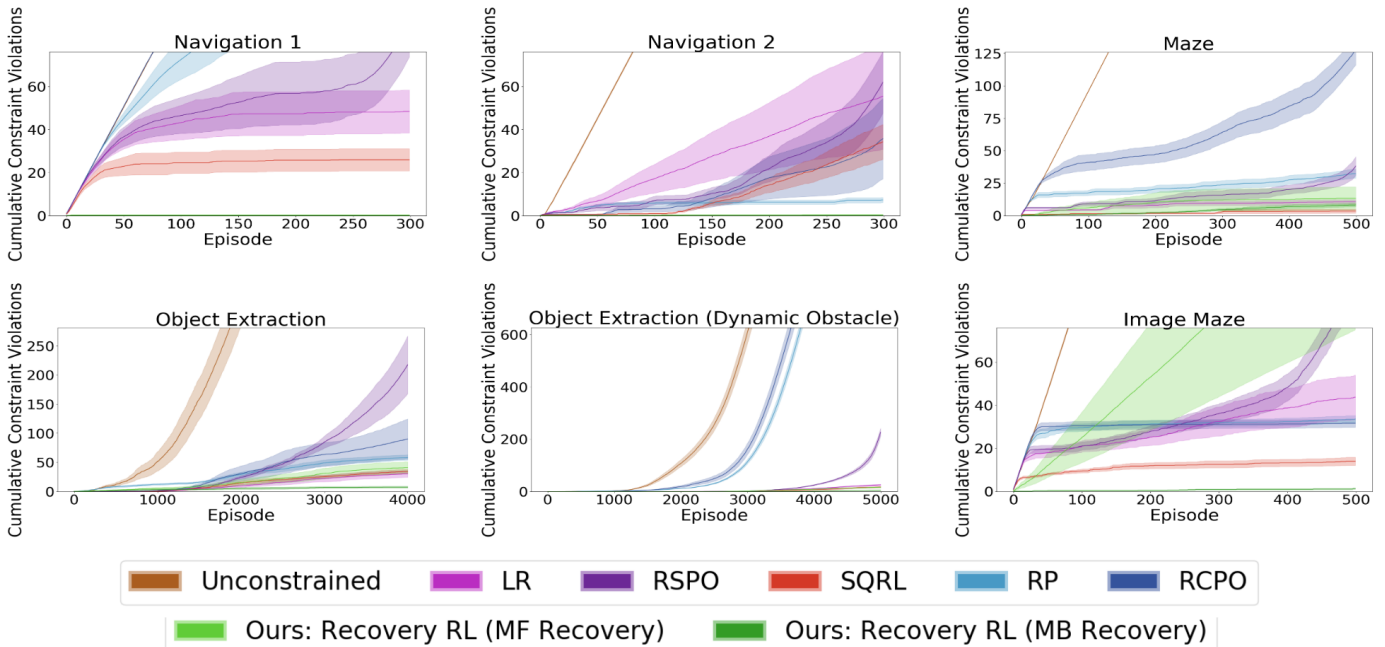


Figure 9: **Simulation Experiments Cumulative Violations:** We plot the cumulative constraint violations for each algorithm in each simulation domain, with results averaged over 10 runs for all algorithms. We observe that Recovery RL (green), is among the safest algorithms across all domains. In the cases where it is less safe than a comparison, it has a higher task success rate (Figure 8).

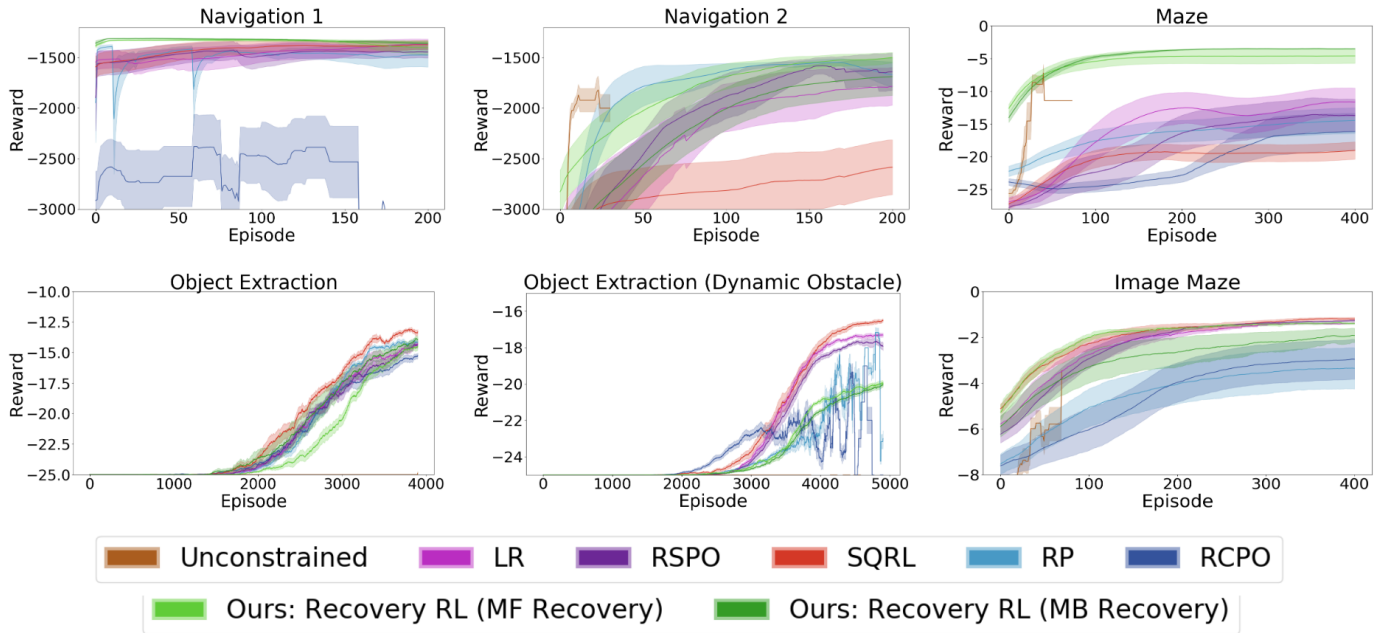


Figure 10: **Simulation Experiments Reward Learning Curve:** We show the total reward attained in each episode smoothed over a 100 episode length window for each simulation domain, with results averaged over 10 runs for all algorithms. We do not show results when constraints are violated to prevent negative bias for algorithms which may violate constraints very frequently, which accounts for gaps in the plot, especially for the unconstrained algorithm which tends to violate constraints very frequently. Thus, the plots illustrate the attained reward for all algorithms conditioned on not violating constraints, which provides a good measure on the quality of solutions found. We find that on all but the dynamic obstacle domain, Recovery RL is able to converge more quickly to higher quality solutions with respect to the task reward function compared to comparisons. This indicates that Recovery RL is able to learn more efficiently, in addition to more safely, compared to comparison algorithms.

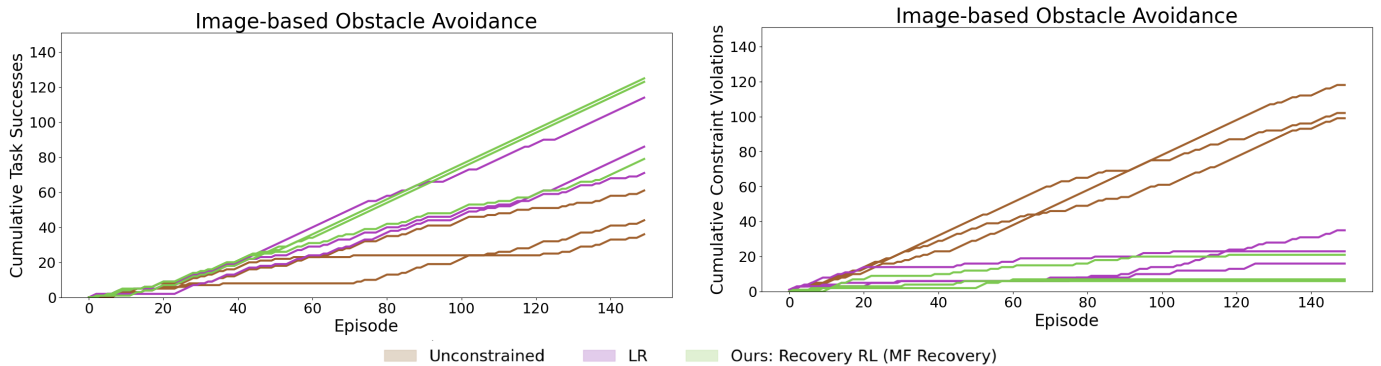


Figure 11: **Physical Experiment Successes and Violations:** We plot the cumulative constraint violations and task successes for the image-based obstacle avoidance task on the dVRK for all 3 runs of each algorithm. We observe that Recovery RL is generally both more successful and safer than LR and unconstrained.

A. Navigation Environments

The Navigation 1 and 2 environments have linear Gaussian dynamics and are built from scratch while the Maze environment is built on the Maze environment from [24]. In all navigation environments, offline data is collected by repeatedly initializing the pointmass agent randomly in the environment and executing controls to make it collide with the nearest obstacle.

- 1) **Navigation 1 and 2:** This environment has single integrator dynamics with additive Gaussian noise sampled from $\mathcal{N}(0, \sigma^2 I_2)$ where $\sigma = 0.05$ and drag coefficient 0.2. The start location is sampled from $\mathcal{N}((-50, 0)^T, I_2)$ and the task is considered successfully completed if the agent gets within 1 unit of the origin. We use negative Euclidean distance from the goal as a reward function. Methods that use a safety critic are given 8000 transitions of data for offline pretraining. For Navigation 1, 455 of these transitions contain constraint violating states, while in Navigation 2, 778 of these transitions contain constraint violating states.
- 2) **Maze:** This environment is implemented in MuJoCo and we again use negative Euclidean distance from the goal as a reward function. Methods that use a safety critic are given 10,000 transitions of data for offline pretraining, 1163 of which contain constraint violating states.

B. Manipulation Environments

We build two manipulation environments on top of the cartgripper environment in the visual foresight repository [9]. The robot can translate in cardinal directions and open/close its grippers. In manipulation environments, offline data is collected by tuning a proportional controller to guide the robot end effector towards the objects. For the offline constraint violations, Gaussian noise is added to the controls when the end effector is sufficiently close to the objects to increase the likelihood of constraint violations. Additionally, to seed the task critic function for SAC to ease exploration for all algorithms, we utilize the same PID controller to collect task demos illustrating the red object being successfully lifted by automatically opening and closing the gripper when the end effector is sufficiently close to the red object.

- 1) **Object Extraction:** This environment is implemented in MuJoCo, and the reward function is -1 until the object is grasped and lifted, at which point it is 0 and the episode terminates. Constraint violations are determined by checking whether any object’s orientation is rotated about the x or y axes by at least 15 degrees. All methods that use a safety critic are given 20,000 transitions of data for offline pretraining, 363 of which contain constraint violating states. All methods are given

Table I: Constraint Violations Breakdown: We report the proportion of constraint violations for each environment that occur when the recovery policy is activated for Recovery RL (format is mean \pm standard error). If most constraint violations occur when the recovery policy is active, this indicates that the safety critic is sufficiently accurate to detect that the recovery policy must be activated, but may not provide sufficient information to avoid constraint violations. We note that if the safety critic detects the need for recovery behavior too late, then these errors are attributed to the recovery policy. For the both Maze environments and the Object Extraction environment, most constraint violations occur when the recovery policy is activated. In Navigation 1, none occur when the recovery policy is activated, but in this environment constraints are almost never violated. In the Image-Based obstacle avoidance tasks, most violations occur when the recovery policy is not activated, which indicates that the bottleneck in this task is the quality of the safety critic. In Navigation 2, Recovery RL never violates constraints and only model-free recovery was run for Recovery RL on the physical robot. In the Dynamic Obstacle Object Extraction environment, we observe a more even combination of low safety critic values and recovery errors during constraint violations.

	Navigation 1	Navigation 2	Maze	Object Extraction	Object Extraction (Dynamic Obstacle)	Image Maze	Image Obstacle Avoidance
MF Recovery	N/A	N/A	0.828 \pm 0.115	0.954 \pm 0.024	0.550 \pm 0.049	0.717 \pm 0.156	0.000 \pm 0.000
MB Recovery	N/A	1.000 \pm 0.000	0.858 \pm 0.039	0.98344 \pm 0.01655	0.269 \pm 0.055	0.583 \pm 0.059	N/A

Table II: Hyperparameters for Recovery RL and comparisons for all domains

	LR	RP	RCPO	MF Recovery	MB Recovery
Navigation 1	(0.8, 0.3, 5000)	1000	(0.8, 0.3, 1000)	(0.8, 0.3)	(0.8, 0.3, 5)
Navigation 2	(0.65, 0.2, 1000)	3000	(0.65, 0.2, 5000)	(0.65, 0.2)	(0.65, 0.2, 5)
Maze	(0.5, 0.15, 100)	50	(0.5, 0.15, 50)	(0.5, 0.15)	(0.5, 0.15, 15)
Object Extraction	(0.75, 0.25, 50)	50	(0.75, 0.25, 50)	(0.75, 0.25)	(0.85, 0.35, 15)
Object Extraction (Dyn. Obstacle)	(0.85, 0.25, 20)	25	(0.85, 0.25, 10)	(0.85, 0.35)	(0.85, 0.25, 15)
Image Maze	(0.65, 0.1, 10)	20	(0.65, 0.1, 20)	(0.65, 0.1)	(0.6, 0.05, 10)
Image Obstacle Avoidance	(0.55, 0.05, 1000)	N/A	N/A	(0.55, 0.05)	N/A

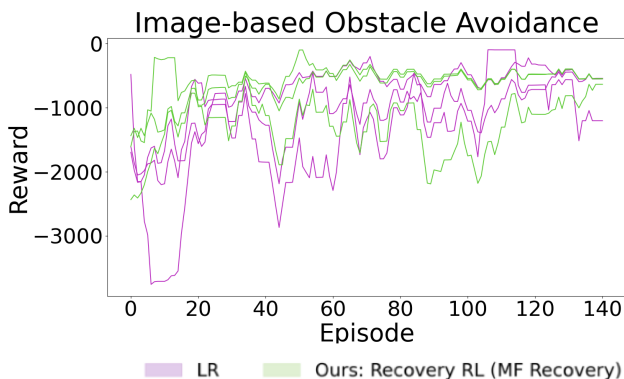


Figure 12: Physical Experiment Reward Learning Curve: We show the total reward attained in each episode smoothed over a 10 episode length window with results from 3 runs for all algorithms. We do not show results when constraints are violated to prevent negative bias for algorithms which may violate constraints very frequently, which accounts for gaps in the plot or explains why some plots have a single line (only one out of 3 runs is constraint satisfying). Note that the Unconstrained algorithm does not appear in the plot as it never makes progress on the harder initial configuration of the task. Thus, the plots illustrate the attained reward for all algorithms conditioned on not violating constraints, which provides a good measure on the quality of solutions found. We find that Recovery RL is able to converge more quickly to higher quality solutions with respect to the task reward function compared to comparisons. This indicates that Recovery RL is able to learn more efficiently, in addition to more safely, compared to comparison algorithms.

Table III: Hyperparameters for Recovery RL and all comparisons.

Algorithm Name	Hyperparameter Format
LR	$(\gamma_{risk}, \epsilon_{risk}, \lambda)$
RP	λ
RCPO	$(\gamma_{risk}, \epsilon_{risk}, \lambda)$
MF Recovery	$(\gamma_{risk}, \epsilon_{risk})$
MB Recovery	$(\gamma_{risk}, \epsilon_{risk}, H)$

1000 transitions of task demonstration data to pretrain the task policy’s critic function.

- Object Extraction (Dynamic Obstacle):** This environment is implemented in MuJoCo, and the reward function is -1 until the object is grasped and lifted, at which point it is 0 and the episode terminates. Constraint violations are determined by checking whether any object’s orientation is rotated about the x or y axes by at least 15 degrees. Additionally, there is a distractor arm that is moving back and forth in the workspace in a periodic fashion. Arm collisions are also considered constraint violations. All methods that use a safety critic are given 20,000 transitions of data for offline pretraining, 896 of which contain constraint violating states. All methods are given 1000 transitions of task demonstration data to pretrain the task policy’s critic function.

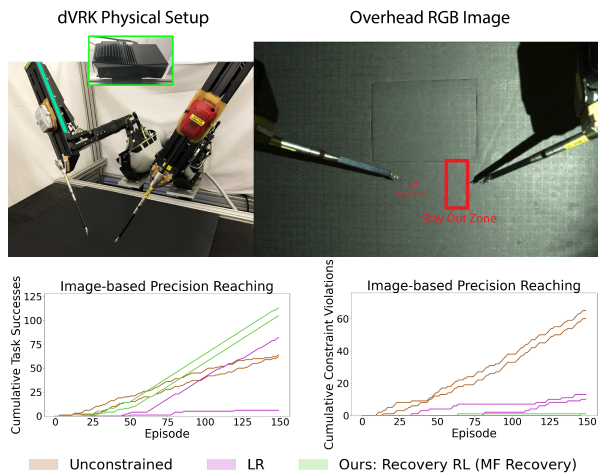


Figure 13: Additional Physical Experiment: The image reacher task on the dVRK involves guiding the end effector to a target position while avoiding an invisible stay out zone in the center of the workspace. We plot the cumulative constraint violations and task successes the image reacher task on the dVRK. We observe that Recovery RL is both more successful and safer than LR and unconstrained.

C. Image Maze

This maze is also implemented in MuJoCo with different walls from the maze that has ground-truth state. Constraint violations occur if the robot collides with a wall. All methods are only supplied with RGB images as input, and all methods that use the safety critic are supplied with 20,000 transitions for pretraining, 3466 of which contain constraint violating states.

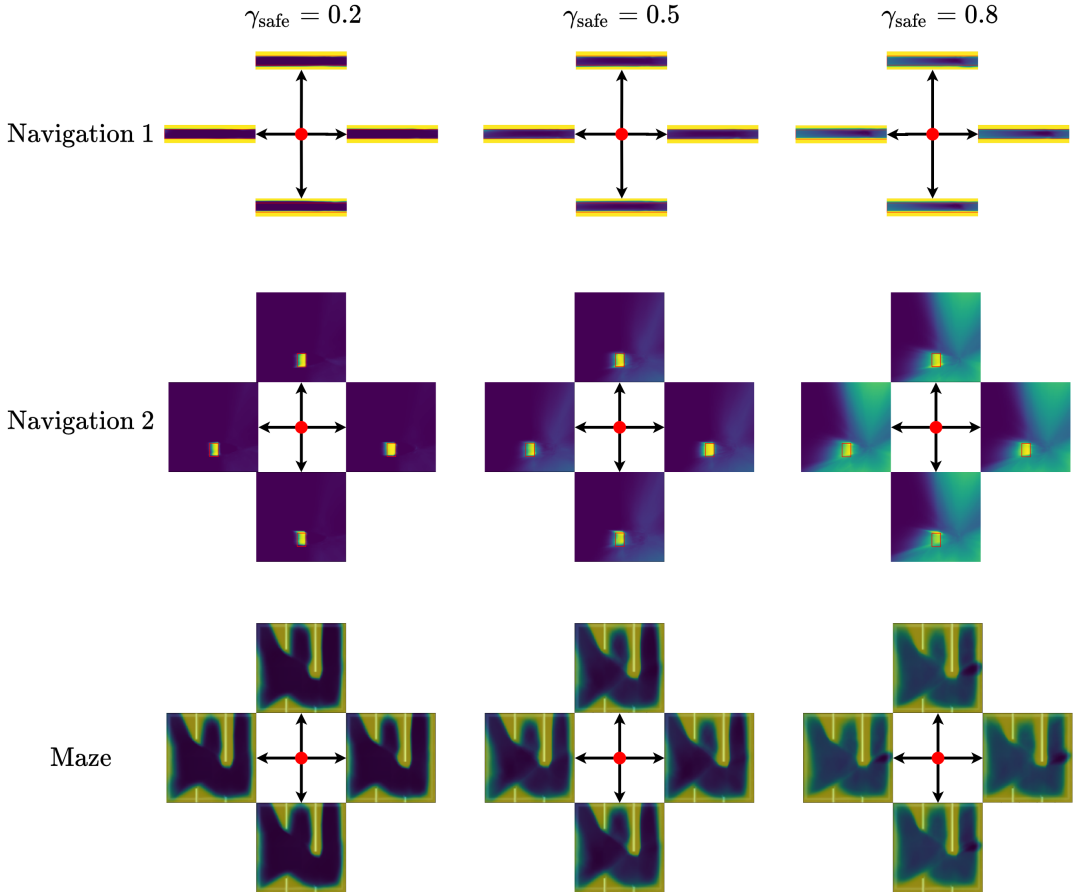


Figure 14: $\hat{Q}_{\phi, \text{risk}}^{\pi}$ **Visualization:** We plot the safety critic Q_{risk}^{π} for the navigation environments using the cardinal directions (left, right, up, down) as action input. We see that as γ_{risk} is increased, the gradient is lower, and the function more gradually increases as it approaches the obstacles. Increasing γ_{risk} essentially increases the amount of information preserved from possible future constraint violations, allowing them to be detected earlier. These plots also illustrate action conditioning of the safety critic values. For example, the down action marks states as more unsafe than the up action directly above walls and obstacles.

D. Physical Experiments

Physical experiments are run on the da Vinci Research Kit (dVRK) [20], a cable-driven bilateral surgical robot. Observations are recorded and supplied to the policies from a Zivid OnePlus RGBD camera. However, we only use RGB images, as the capture rate is much faster, and we subsample the images so input images have dimensions $48 \times 64 \times 3$. End effector position is checked by the environment using the robot’s odometry to check task completion, but this is not supplied to any of the policies. In practice, the robot’s end effector position can be slightly inaccurate due to cabling effects such as hysteresis [18], but we ignore these effects in this paper. We train a neural network to classify whether the robot is in collision based on its current readings and joint position. All methods that use a safety critic are supplied with 6,000 transitions of data for pretraining, 649 of which contain constraint violating states. As for the navigation environments, offline data is collected by randomly initializing the end effector in the environment and guiding it

towards the nearest obstacle. To reduce extrapolation errors during learning, we sample a start state on the right and left sides of the workspace with equal probability.

E. Additional Physical Experiment

We also evaluate Recovery RL and comparisons on an image-based reaching task where the robot must make sure the end effector position does not intersect with a stay out zone in the center of the workspace instead of physical bumpers. The setup is almost identical to the setup described in Section XIV-D. We again provide RGB images to algorithms, and use 10,000 transitions to pre-train the safety critic. We again find that Recovery RL is able to outperform comparisons on this task, both in terms of constraint satisfaction, and task completion.