

# Safe Learning in Robotics: From Learning-Based Control to Safe Reinforcement Learning

Lukas Brunke\*, Melissa Greeff\*, Adam W. Hall\*, Zhaocong Yuan\*, Siqi Zhou\*, Jacopo Panerati, and Angela P. Schoellig

\*Equal contribution

Institute for Aerospace Studies, University of Toronto, Toronto, Ontario, Canada, M3H 5T6; University of Toronto Robotics Institute, Toronto, Ontario, Canada, M5S 1A4; Vector Institute for Artificial Intelligence, Toronto, Ontario, Canada, M5G 1M1; emails: {firstname.lastname}@robotics.utias.utoronto.ca

Published in the Annual Review of Control, Robotics, and Autonomous Systems.

This unofficial version compiled on: December 8, 2021.

## Keywords

safe learning, robotics, robot learning, learning-based control, safe reinforcement learning, adaptive control, robust control, model predictive control, machine learning, benchmarks

## Abstract

The last half-decade has seen a steep rise in the number of contributions on safe learning methods for real-world robotic deployments from both the control and reinforcement learning communities. This article provides a concise but holistic review of the recent advances made in using machine learning to achieve safe decision making under uncertainties, with a focus on unifying the language and frameworks used in control theory and reinforcement learning research. Our review includes: learning-based control approaches that safely improve performance by learning the uncertain dynamics, reinforcement learning approaches that encourage safety or robustness, and methods that can formally certify the safety of a learned control policy. As data- and learning-based robot control methods continue to gain traction, researchers must understand when and how to best leverage them in real-world scenarios where safety is imperative, such as when operating in close proximity to humans. We highlight some of the open challenges that will drive the field of robot learning in the coming years, and emphasize the need for realistic physics-based benchmarks to facilitate fair comparisons between control and reinforcement learning approaches.

## Contents

1. INTRODUCTION .....	2
2. PRELIMINARIES AND BACKGROUND OF SAFE LEARNING CONTROL .....	3
2.1. Problem Statement .....	4
2.2. A Control Theory Perspective .....	6
2.3. A Reinforcement Learning Perspective .....	8
2.4. Bridging Control Theory and RL for Safe Learning Control .....	9
3. SAFE LEARNING CONTROL APPROACHES .....	10
3.1. Learning Uncertain Dynamics to Safely Improve Performance .....	11
3.2. Encouraging Safety and Robustness in Reinforcement Learning .....	14
3.3. Certifying Learning-Based Control Under Dynamics Uncertainty .....	18
4. BENCHMARKS .....	22
5. DISCUSSION AND PERSPECTIVES FOR FUTURE DIRECTIONS .....	25
A. SUMMARY OF REVIEWED LITERATURE .....	34
B. REVISION HISTORY .....	36

## 1. INTRODUCTION

Robotics researchers strive to design systems that can operate autonomously in increasingly complex scenarios, often in close proximity to humans. Examples include self-driving vehicles (1), aerial delivery (2), and the use of mobile manipulators for service tasks (3). However, the dynamics of these complex applications are often uncertain or only partially known—for example, the mass distribution of a carried payload might not be given *a priori*. Uncertainties arise from various sources. For example, the robot dynamics may not be perfectly modeled, sensor measurements may be noisy, and/or the operating environment may not be well-characterized or may include other agents whose dynamics and plans are not known.

In these real-world applications, robots must *make decisions despite having only partial knowledge of the world*. In recent years, the research community has multiplied its efforts to leverage data-based approaches to address this problem. This was motivated in part by the success of machine learning in other areas such as computer vision and natural language processing.

A crucial, domain-specific challenge of *learning for robot control* is the need to implement and formally guarantee safety of the robot’s behavior, not only for the optimized policy (or controller, which is essential for the certification of systems that interact with humans) but also during learning, to avoid costly hardware failures and improve convergence. Ultimately, these safety guarantees can only be derived from the assumptions and structure captured by the problem formalization.

Both control theory and machine learning—reinforcement learning (RL) in particular—have proposed approaches to tackle this problem. Control theory has traditionally taken a model-driven approach (see **Figure 1**): it leverages a given dynamics model and provides guarantees with respect to known operating conditions. RL has traditionally taken a data-driven approach, which makes it highly adaptable to new contexts at the expense of providing formal guarantees. Combining model-driven and data-driven approaches, and leveraging the advantages of each, is a promising direction for safe learning in robotics. The methods we review encourage **robustness** (by accounting for the worst-case scenarios and taking conservative actions), enable **adaptation** (by learning from online observations and adapting to unknown situations), and build and leverage **prediction models** (based on a combination of domain knowledge, real-world data, and high-fidelity simulators).

While control is still the bedrock of all current robot applications, the body of safe RL

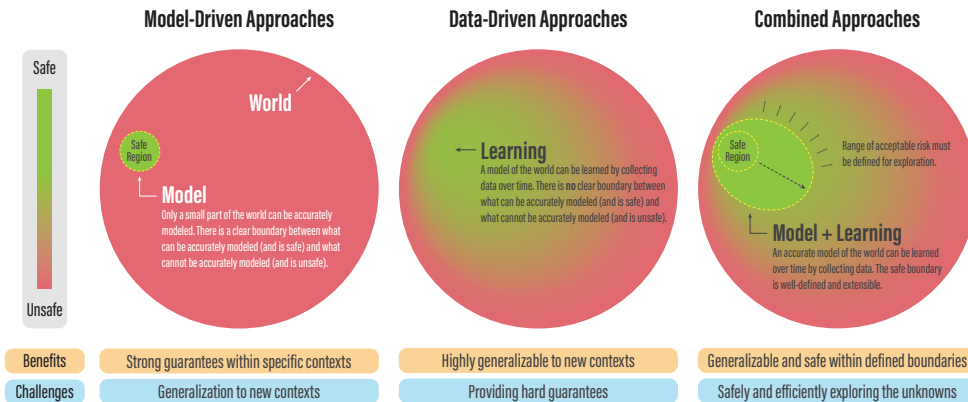


Figure 1: A comparison of model-driven, data-driven, and combined approaches.

literature has ballooned from tens to over a thousand publications in just a few years since its most recent review (4). Physics-based simulation (5), which we leverage in our open-source benchmark implementation<sup>1</sup>, has played an important role in the recent progress of RL, however, the transfer to real systems remains a research area in itself (6).

Previous review works focused on specific techniques—for example, learning-based model predictive control (MPC) (7), iterative learning control (ILC) (8, 9), model-based RL (10), data-efficient policy search (11), imitation learning (12), or the use of RL in robotics (13, 14) and in optimal control (15)—without emphasizing the safety aspect. Recent surveys on safe learning control focus on either control-theoretic (16) or RL approaches (17), and do not provide a unifying perspective.

In this article, we provide a bird’s eye view of the most recent work in learning-based control and reinforcement learning that implement safety and provide safety guarantees for robot control. We focus on safe learning control approaches where the data generated by the robot system is used to learn or modify the feedback controller (or policy). We hope to help shrink the gap between the control and RL communities by creating a common vocabulary and introducing benchmarks for algorithm evaluation that can be leveraged by both (18, 19). Our target audience are researchers, with either a control or RL background, who are interested in a concise but holistic perspective on the problem of safe learning control. While we do not cover perception, estimation, planning, or multi-agent systems, we do connect our discussion to these additional challenges and opportunities.

## 2. PRELIMINARIES AND BACKGROUND OF SAFE LEARNING CONTROL

In this review, we are interested in the problem of *safe decision making under uncertainties using machine learning* (i.e., *safe learning control*). Intuitively, in safe learning control, our goal is to allow a robot to fulfil a task while respecting a set of safety constraints despite the uncertainties present in the problem. In this section, we define the safe learning control problem (Sec. 2.1) and provide an overview of how the problem of *decision making under uncertainties* has traditionally been tackled by the control (Sec. 2.2) and RL communities (Sec. 2.3). We highlight the main limitations of these approaches and articulate how novel data-based, safety-focused methods can address these gaps (Sec. 2.4).

<sup>1</sup>Safe control benchmark suite on GitHub: <https://github.com/utiasDSL/safe-control-gym>

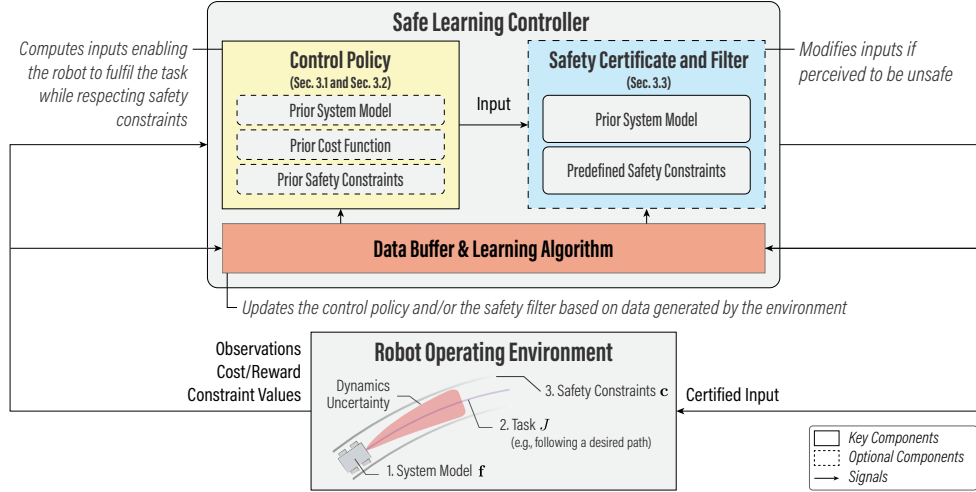


Figure 2: The safe learning control problem is defined by the cost function  $J$ , the system model  $\mathbf{f}$ , and the constraints  $\mathbf{c}$ , which may all be initially unknown. Data is used to update the control policy (see Sec. 3.1 and Sec. 3.2) or the safety filter (see Sec. 3.3).

### M1. Transition

#### Probability Function:

The dynamics model in Equation 1 can be equivalently represented as a state transition probability function  $T_k(\mathbf{x}_{k+1} | \mathbf{x}_k, \mathbf{u}_k)$ , commonly seen in RL approaches (20).

## 2.1. Problem Statement

We formulate the safe learning control problem as an optimization problem to capture the efforts of both the control and RL communities. The optimization problem has three main components (see Figure 2): (i) a system model that describes the dynamic behavior of the robot, (ii) a cost function that defines the control objective or task goal, and (iii) a set of constraints that specify the safety requirements. The goal is to find a *controller* or *policy* that computes commands (also called *inputs*) that enable the system to fulfil the task while respecting given safety constraints. In general, any of the components could be initially unknown or only partially known. Below, we first introduce each of the three components and conclude by stating the overall safe learning control problem.

**2.1.1. System Model.** We consider a robot whose dynamics can be represented by the following discrete-time model:

$$\mathbf{x}_{k+1} = \mathbf{f}_k(\mathbf{x}_k, \mathbf{u}_k, \mathbf{w}_k), \quad 1.$$

where  $k \in \mathbb{Z}_{\geq 0}$  is the discrete-time index,  $\mathbf{x}_k \in \mathbb{X}$  is the state with  $\mathbb{X}$  denoting the state space,  $\mathbf{u}_k \in \mathbb{U}$  is the input with  $\mathbb{U}$  denoting the input space,  $\mathbf{f}_k$  is the dynamics model of the robot,  $\mathbf{w}_k \in \mathbb{W}$  is the process noise distributed according to a distribution  $\mathcal{W}$  (see also M1). Throughout this review, we assume direct access to (possibly noisy) measurements of the robot state  $\mathbf{x}_k$  and neglect the problem of state estimation. Equation 1 represents many common robot platforms (e.g., quadrotors, manipulators, and ground vehicles). More complex models (e.g., partial differential equations) may be necessary for other robot designs.

**2.1.2. Cost Function.** The robot's task is defined by a cost function. We consider a finite-horizon optimal control problem with time horizon  $N$ . Given an initial state  $\mathbf{x}_0$ , the cost is computed based on the sequence of states  $\mathbf{x}_{0:N} = \{\mathbf{x}_0, \mathbf{x}_1, \dots, \mathbf{x}_N\}$  and the sequence of inputs  $\mathbf{u}_{0:N-1} = \{\mathbf{u}_0, \mathbf{u}_1, \dots, \mathbf{u}_{N-1}\}$ :

$$J(\mathbf{x}_{0:N}, \mathbf{u}_{0:N-1}) = l_N(\mathbf{x}_N) + \sum_{k=0}^{N-1} l_k(\mathbf{x}_k, \mathbf{u}_k), \quad 2.$$

where  $l_k : \mathbb{X} \times \mathbb{U} \mapsto \mathbb{R}$  is the stage cost incurred at each time step  $k$ , and  $l_N : \mathbb{X} \mapsto \mathbb{R}$  is the terminal cost incurred at the end of the  $N$ -step horizon (see also M2). The stage and terminal cost functions map the state and input sequences, which may be random variables, to a real number, and may, for example, include the expected value or variance operators.

**2.1.3. Safety Constraints.** Safety constraints ensure, or encourage, the safe operation of the robot and include: (i) *state constraints*  $\mathbb{X}_c \subseteq \mathbb{X}$ , which define the set of safe operating states (e.g., the lane in self-driving), (ii) *input constraints*  $\mathbb{U}_c \subseteq \mathbb{U}$  (e.g., actuation limits), and (iii) *stability guarantees* (e.g., the robot’s motion converging to a desired path, see M3).

To encode the safety constraints, we define  $n_c$  constraint functions:  $\mathbf{c}_k(\mathbf{x}_k, \mathbf{u}_k, \mathbf{w}_k) \in \mathbb{R}^{n_c}$  with each constraint  $c_k^j$  being a real-valued, time-varying function. Starting with the strongest guarantee, we introduce three levels of safety: hard, probabilistic, and soft constraints (illustrated in Figure 3). In practice, safety levels are often mixed. For example, input constraints are typically hard constraints but state constraints may be soft constraints.

**Safety Level III: Constraint Satisfaction Guaranteed.** The system satisfies hard constraints:

$$c_k^j(\mathbf{x}_k, \mathbf{u}_k, \mathbf{w}_k) \leq 0, \quad 3.$$

for all times  $k \in \{0, \dots, N\}$  and constraint indexes  $j \in \{1, \dots, n_c\}$ .

**Safety Level II: Constraint Satisfaction with Probability  $p$ .** The system satisfies probabilistic constraints:

$$\Pr \left( c_k^j(\mathbf{x}_k, \mathbf{u}_k, \mathbf{w}_k) \leq 0 \right) \geq p^j, \quad 4.$$

where  $\Pr(\cdot)$  denotes the probability and  $p^j \in (0, 1)$  defines the likelihood of the  $j$ -th constraint being satisfied, with  $j \in \{1, \dots, n_c\}$  and for all times  $k \in \{0, \dots, N\}$ . The chance constraint in Equation 4 is identical to the hard constraint in Equation 3 for  $p^j = 1$ .

**Safety Level I: Constraint Satisfaction Encouraged.** The system encourages constraint satisfaction. This can be achieved in different ways. One way is to add a penalty term to the objective function that discourages the violation of constraints with a high cost. A non-negative  $\epsilon_j$  is added to the right-hand side of the inequality Equation 3, for all times  $k \in \{0, \dots, N\}$  and  $j \in \{1, \dots, n_c\}$ ,

$$c_k^j(\mathbf{x}_k, \mathbf{u}_k, \mathbf{w}_k) \leq \epsilon_j, \quad 5.$$

and an appropriate penalty term  $l_\epsilon(\epsilon) \geq 0$  with  $l_\epsilon(\epsilon) = 0 \iff \epsilon = \mathbf{0}$  is added to the objective function. The vector  $\epsilon$  includes all elements  $\epsilon_j$  and is an additional variable of the optimization problem. Alternatively, although  $c_k^j(\mathbf{x}_k, \mathbf{u}_k, \mathbf{w}_k)$  is a step-wise quantity, some approaches only aim to provide guarantees on its expected value  $E[\cdot]$  on a trajectory level:

$$J_{c^j} = E \left[ \sum_{k=0}^{N-1} c_k^j(\mathbf{x}_k, \mathbf{u}_k, \mathbf{w}_k) \right] \leq d_j, \quad 6.$$

where  $J_{c^j}$  represents the expected total constraint cost, and  $d_j$  defines the corresponding constraint threshold. The constraint function can optionally be discounted as  $\gamma^k c_k^j(\mathbf{x}_k, \mathbf{u}_k, \mathbf{w}_k)$ , similar to the stage cost (see M2).

**2.1.4. Formulation of the Safe Learning Control Problem.** The functions introduced above, i.e., the system model  $\mathbf{f}$ , the constraints  $\mathbf{c}$ , and the cost function  $J$ , represent the *true* functions of the robot control problem. In practice,  $\mathbf{f}$ ,  $\mathbf{c}$ , and  $J$  may be unknown or partially

---

**M2. Discounted Rewards:** The stage cost can be related to the discounted rewards in RL:  $l_k = -\gamma^k r_k(\mathbf{x}_k, \mathbf{u}_k)$ , where  $r_k : \mathbb{X} \times \mathbb{U} \mapsto \mathbb{R}$  is the reward function, and  $\gamma \in [0, 1]$  is the discount factor (20). We formulate a cost minimization problem for safe learning control, while RL typically solves a reward maximization problem.

---

**M3. Stability:** There are different notions of stability in the control literature (21). Stability generally implies boundedness of the system state. One common notion of stability is Lyapunov stability, which requires the system to either stay close to (*stable*) or converge to (*asymptotically stable*) a desired state.

---

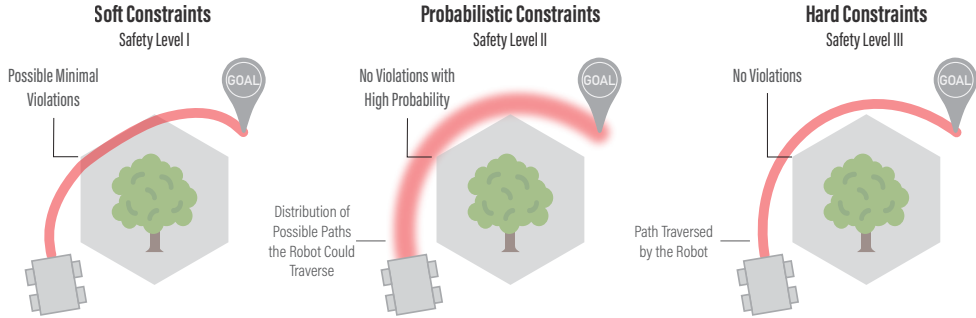


Figure 3: Illustration of the different safety levels.

known. Without loss of generality, we assume that each of the true functions  $\mathbf{f}$ ,  $\mathbf{c}$ , and  $J$  can be decomposed into a nominal component  $(\bar{\cdot})$ , reflecting our prior knowledge, and an unknown component  $(\hat{\cdot})$ , to be learned from data. For instance, the dynamics model  $\mathbf{f}$  can be decomposed as

$$\mathbf{f}_k(\mathbf{x}_k, \mathbf{u}_k, \mathbf{w}_k) = \bar{\mathbf{f}}_k(\mathbf{x}_k, \mathbf{u}_k) + \hat{\mathbf{f}}_k(\mathbf{x}_k, \mathbf{u}_k, \mathbf{w}_k), \quad 7.$$

where  $\bar{\mathbf{f}}$  is the prior dynamics model and  $\hat{\mathbf{f}}$  are the uncertain dynamics.

*Safe learning control* (SLC) leverages our prior knowledge  $\mathcal{P} = \{\bar{\mathbf{f}}, \bar{\mathbf{c}}, \bar{J}\}$  and the data collected from the system  $\mathcal{D} = \{\mathbf{x}^{(i)}, \mathbf{u}^{(i)}, \mathbf{c}^{(i)}, l^{(i)}\}_{i=0}^D$  to find a policy (or controller)  $\pi_k(\mathbf{x}_k)$  that achieves the given task while respecting all safety constraints,

$$\text{SLC} : (\mathcal{P}, \mathcal{D}) \mapsto \pi_k, \quad 8.$$

where  $(\cdot)^{(i)}$  denotes a sample of a quantity  $(\cdot)_k$  and  $D$  is the data set size. More specifically, we aim to find a policy  $\pi_k$  that best approximates the true optimal policy  $\pi_k^*$ , which is the solution to the following optimization problem:

$$J^{\pi^*}(\bar{\mathbf{x}}_0) = \min_{\pi_{0:N-1}, \epsilon} J(\mathbf{x}_{0:N}, \mathbf{u}_{0:N-1}) + l_\epsilon(\epsilon) \quad 9a.$$

$$\text{s.t. } \mathbf{x}_{k+1} = \mathbf{f}_k(\mathbf{x}_k, \mathbf{u}_k, \mathbf{w}_k), \quad \mathbf{w}_k \sim \mathcal{W}, \quad \forall k \in \{0, \dots, N-1\}, \quad 9b.$$

$$\mathbf{u}_k = \pi_k(\mathbf{x}_k), \quad 9c.$$

$$\mathbf{x}_0 = \bar{\mathbf{x}}_0, \quad 9d.$$

$$\text{Safety constraints according to either} \quad 9e.$$

Equation 3, Equation 4, Equation 5 or Equation 6, and  $\epsilon \geq 0$ ,

where  $\bar{\mathbf{x}}_0 \sim \mathcal{X}_0$  is the initial state with  $\mathcal{X}_0$  being the initial state distribution, and  $\epsilon$  and  $l_\epsilon$  are introduced to account for the soft safety constraint case (Safety Level I, Equation 5) and are set to zero, for example, if only hard and probabilistic safety constraints are considered (Safety Levels II and III).

## 2.2. A Control Theory Perspective

Safe decision making under uncertainty has a long history in the field of control theory. Typical assumptions are that a model of the system is available and it is either parameterized by an unknown parameter or it has bounded unknown dynamics and noise. While the control approaches are commonly formulated using continuous-time dynamic models, they are usually implemented in discrete time with sampled inputs and measurements (see M4).

**Adaptive control** considers systems with parametric uncertainties (see M5) and adapts the controller parameters online to optimize performance. Adaptive control requires knowledge of the parametric form of the uncertainty (23) and typically considers a dynamics

### M4. Continuous Time:

Many classical control approaches rely on a continuous-time formulation. Consider the uncertain time-varying, continuous-time model  $\dot{\mathbf{x}} = \mathbf{f}_t(\mathbf{x}, \mathbf{u}, \mathbf{w})$ , where  $\mathbf{x}$ ,  $\mathbf{u}$  and  $\mathbf{w}$  are functions of time  $t$ . We can recover our discrete state through sampling  $\mathbf{x}_k = \mathbf{x}(k\Delta t)$ , where  $\Delta t$  is the sampling time. Similarly, the discrete control input is  $\mathbf{u}_k = \mathbf{u}(k\Delta t)$  and often kept constant (22) over the time interval  $[k\Delta t, (k+1)\Delta t)$ .

model that is affine in  $\mathbf{u}$  and the uncertain parameters  $\boldsymbol{\theta} \in \Theta$ :

$$\mathbf{x}_{k+1} = \bar{\mathbf{f}}_{\mathbf{x}}(\mathbf{x}_k) + \bar{\mathbf{f}}_{\mathbf{u}}(\mathbf{x}_k)\mathbf{u}_k + \bar{\mathbf{f}}_{\boldsymbol{\theta}}(\mathbf{x}_k)\boldsymbol{\theta}, \quad 10.$$

where  $\bar{\mathbf{f}}_{\mathbf{x}}$ ,  $\bar{\mathbf{f}}_{\mathbf{u}}$ , and  $\bar{\mathbf{f}}_{\boldsymbol{\theta}}$  are known functions often derived from first principles and  $\Theta$  is a possibly bounded parameter set. The control input is  $\mathbf{u}_k = \boldsymbol{\pi}(\mathbf{x}_k, \hat{\boldsymbol{\theta}}_k)$ , which is parameterized by  $\hat{\boldsymbol{\theta}}_k$ . The parameter  $\hat{\boldsymbol{\theta}}_k$  is adapted by using either a *Lyapunov function* to guarantee that the closed-loop system is stable (see M6) or *Model Reference Adaptive Control (MRAC)* to make the system behave as a predefined stable reference model (23). Adaptive control is typically limited to parametric uncertainties and relies on a specific model structure. Moreover, adaptive control approaches tend to “overfit” to the latest observations and convergence to the true parameters is generally not guaranteed (23, 24). These limitations motivate the learning-based adaptive control approaches in Sec. 3.1.1.

**Robust control** is a control design technique that guarantees stability for pre-specified bounded disturbances, which can include unknown dynamics and noise. In contrast to adaptive control, which adapts to the parameters currently present, robust control finds a suitable controller for all possible disturbances and keeps the controller unchanged after the initial design. Robust control is largely limited to linear, time-invariant (LTI) systems with linear nominal model  $\bar{\mathbf{f}}(\mathbf{x}_k, \mathbf{u}_k) = \bar{\mathbf{A}}\mathbf{x}_k + \bar{\mathbf{B}}\mathbf{u}_k$  and unknown dynamics  $\hat{\mathbf{f}}_k(\mathbf{x}_k, \mathbf{u}_k, \mathbf{w}_k) = \hat{\mathbf{A}}\mathbf{x}_k + \hat{\mathbf{B}}\mathbf{u}_k + \mathbf{w}_k \in \mathbb{D}$ , with  $\mathbb{D}$  being known and bounded, and  $\bar{\mathbf{A}}$ ,  $\bar{\mathbf{B}}$ ,  $\hat{\mathbf{A}}$ ,  $\hat{\mathbf{B}}$  being static matrices of appropriate size. That is,

$$\mathbf{x}_{k+1} = (\bar{\mathbf{A}} + \hat{\mathbf{A}})\mathbf{x}_k + (\bar{\mathbf{B}} + \hat{\mathbf{B}})\mathbf{u}_k + \mathbf{w}_k. \quad 11.$$

Robust control design techniques, such as robust  $H_{\infty}$ - and  $H_2$ -control design (25), yield controllers that are robustly stable for all  $\hat{\mathbf{f}}_k \in \mathbb{D}$ . Robust control can be extended to nonlinear systems whose dynamics can be decomposed into a linear nominal model  $\bar{\mathbf{f}}$  and a nonlinear function  $\hat{\mathbf{f}}$  with known bound  $\hat{\mathbf{f}} \in \mathbb{D}$  (26).

**Robust Model Predictive Control (MPC)** extends classical adaptive and robust control by additionally guaranteeing state and input constraints,  $\mathbf{x}_k \in \mathbb{X}_c$  and  $\mathbf{u}_k \in \mathbb{U}_c$ , for all possible bounded disturbances  $\hat{\mathbf{f}} \in \mathbb{D}$ . At every time step  $k$ , MPC solves a constrained optimization problem over a control input sequence  $\mathbf{u}_{k:k+H-1}$  for a finite horizon  $H$ , and applies the first optimal control input to the system and then re-solves the optimization problem in the next time step based on the current state (27). A common approach in robust MPC (RMPC) is *tube-based MPC* (28), which uses a nominal prediction model  $\bar{\mathbf{f}}(\mathbf{x}_k, \mathbf{u}_k)$  in the MPC optimization and tightens the constraints to account for unmodeled dynamics. A stabilizing controller keeps the true state inside a bounded set of states around the nominal state, called *tube*, for all possible disturbances. Since the nominal states satisfy the tightened constraints and the true states stay inside the tube around the nominal states, constraint satisfaction for the true states is guaranteed. Tube-based MPC typically considers a linear nominal model  $\bar{\mathbf{f}}(\bar{\mathbf{x}}_k, \bar{\mathbf{u}}_k) = \bar{\mathbf{A}}\bar{\mathbf{x}}_k + \bar{\mathbf{B}}\bar{\mathbf{u}}_k$  with nominal state  $\bar{\mathbf{x}}_k$  and input  $\bar{\mathbf{u}}_k$ . In its simplest implementation, prior knowledge of set  $\mathbb{D}$  is combined with a known stabilizing linear controller,  $\mathbf{u}_{k,\text{stab}} = \mathbf{K}(\mathbf{x}_k - \bar{\mathbf{x}}_k)$  with gain  $\mathbf{K}$ , to determine the bounded tube  $\Omega_{\text{tube}}$  from the matrices  $\bar{\mathbf{A}}$ ,  $\bar{\mathbf{B}}$ ,  $\mathbf{K}$ , and the set  $\mathbb{D}$ . The stabilizing controller  $\mathbf{u}_{k,\text{stab}}$  keeps all potential errors within the tube,  $\mathbf{x}_k - \bar{\mathbf{x}}_k \in \Omega_{\text{tube}}$ , for all  $k$ . For the nominal model, tube-based MPC solves the following constrained optimization problem at every time step  $k$  to obtain the optimal sequence  $\bar{\mathbf{u}}_{0:H-1}^*$ :

$$J_{\text{RMPC}}^*(\bar{\mathbf{x}}_k) = \min_{\bar{\mathbf{u}}_{0:H-1}} l_H(\mathbf{z}_H) + \sum_{i=0}^{H-1} l_i(\mathbf{z}_i, \bar{\mathbf{u}}_i) \quad 12a.$$

$$\text{s.t. } \mathbf{z}_{i+1} = \bar{\mathbf{A}}\mathbf{z}_i + \bar{\mathbf{B}}\bar{\mathbf{u}}_i, \quad \forall i \in \{0, \dots, H-1\} \quad 12b.$$

$$\mathbf{z}_i \in \mathbb{X}_c \ominus \Omega_{\text{tube}}, \quad \bar{\mathbf{u}}_i \in \mathbb{U}_c \ominus \mathbf{K}\Omega_{\text{tube}}, \quad 12c.$$

$$\mathbf{z}_0 = \bar{\mathbf{x}}_k, \quad \mathbf{z}_H \in \mathbb{X}_{\text{term}}, \quad 12d.$$

#### M5. Parametric Uncertainty: A

*parametric model* depends on a finite number of parameters that may have a physical interpretation or reflect our prior knowledge about the system structure in other ways. *Parametric uncertainty* is the uncertainty in the model parameters.

#### M6. Lyapunov

**Function:** Lyapunov functions are used to analyze the stability of a dynamical system (21). Consider a

closed-loop system under some policy  $\boldsymbol{\pi}(\mathbf{x})$ :

$$\mathbf{x}_{k+1} = \mathbf{f}_{\boldsymbol{\pi}}(\mathbf{x}_k) = \mathbf{f}(\mathbf{x}_k, \boldsymbol{\pi}(\mathbf{x}_k))$$

with  $\mathbf{f}_{\boldsymbol{\pi}}(\mathbf{x}_k)$  being Lipschitz continuous. A Lipschitz continuous positive definite function

$$L: \mathbb{X} \rightarrow \mathbb{R}_{\geq 0} \text{ with } L(\mathbf{0}) = 0 \text{ and}$$

$$L(\mathbf{x}) > 0, \forall \mathbf{x} \neq \mathbf{0},$$

is a Lyapunov function, if  $L$  maps states under closed-loop state feedback to strictly smaller values (i.e.,  $\Delta L(\mathbf{x}) = L(\mathbf{f}_{\boldsymbol{\pi}}(\mathbf{x})) - L(\mathbf{x}) < 0$ ). This implies that the state converges to the equilibrium at the origin. There exists an analogous formulation for continuous-time systems based on the derivatives of  $L$  (21).

---

**M7. Constraint**

**Tightening:** The constraint sets  $\mathbb{X}_c$  and  $\mathbb{U}_c$  are tightened by an error tube  $\Omega_{\text{tube}}$  using the Pontryagin difference,  $\mathbb{X}_c \ominus \Omega_{\text{tube}} = \{\mathbf{x} \in \mathbb{X} : \mathbf{x} + \boldsymbol{\omega} \in \mathbb{X}_c, \forall \boldsymbol{\omega} \in \Omega_{\text{tube}}\}$  and  $\mathbb{U}_c \ominus \mathbf{K}\Omega_{\text{tube}}$ , where  $\mathbf{K}$  maps all elements in  $\Omega_{\text{tube}}$  to the control input space  $\mathbb{U}$  with  $\mathbf{K}\Omega_{\text{tube}}$ .

---

---

**M8. Value and Action-Value**

**Functions:** The value function of state  $\mathbf{x}_k$  under policy  $\pi_k$ ,  $V^{\pi_k}(\mathbf{x}_k)$ , is the expected return  $J$  of applying  $\pi_k$  from  $\mathbf{x}_k$ . Similarly, the action-value function of action  $\mathbf{u}_k$  in state  $\mathbf{x}_k$  under  $\pi_k$ ,  $Q^{\pi_k}(\mathbf{x}_k, \mathbf{u}_k)$ , is the expected return  $J$  when taking the action  $\mathbf{u}_k$  at  $\mathbf{x}_k$ , and then following  $\pi_k$ .

---

---

**M9. Neural**

**Network:** A neural network (NN) is a computational model with interconnected layers of neurons (parameterized by weights) that can be used to approximate highly nonlinear functions.

---

where  $\mathbf{z}_i$  is the open-loop nominal state at time step  $k+i$  and the state and input constraints are tightened using the bounded tube  $\Omega_{\text{tube}}$  (see M7). Combined with the stabilizing control input  $\mathbf{u}_{k,\text{stab}}$ , the control input  $\mathbf{u}_k = \mathbf{u}_{k,\text{stab}} + \bar{\mathbf{u}}_0^*$  is applied to the system at every time step. Stability and satisfaction of the tightened constraints is guaranteed by selecting the terminal cost  $l_H$  in Equation 12a, such that after the prediction horizon  $H$  the nominal state  $\mathbf{z}_H$  is within a terminal constraint set  $\mathbb{X}_{\text{term}}$  (see Equation 12d) at which point a known linear controller can be safely applied (27).

Both robust control and robust MPC are conservative, as they guarantee stability and—in the case of MPC—state and input constraints for the worst-case scenario. This usually yields poor performance, as described in (7). For example, a conservative uncertainty set  $\mathbb{D}$  generates a large tube  $\Omega_{\text{tube}}$  resulting in tight hard constraints that are prioritized over cost optimization. Learning-based robust control and RMPC improve performance by using data to (i) learn a less conservative state- and input-dependent uncertainty set  $\mathbb{D}$  and/or (ii) learn the unknown dynamics  $\hat{\mathbf{f}}$  and, as a result, reduce the remaining model uncertainty, see Sec. 3.1.2 and Sec. 3.1.3 respectively.

### 2.3. A Reinforcement Learning Perspective

Reinforcement learning (RL) is the standard machine learning framework to address the problem of sequential decision making under uncertainty. Unlike traditional control, RL generally does not rely on an *a priori* dynamics model  $\hat{\mathbf{f}}$  and can be directly applied to uncertain dynamics  $\mathbf{f}$ . However, the lack of explicit assumptions and constraints in many of the works limit their applicability to safe control. RL algorithms attempt to find  $\pi^*$  while gathering data and knowledge of  $\mathbf{f}$  from interaction with the system—taking random actions, initially, and improving afterwards. A long-standing challenge of RL, which hampers safety during the learning stages, is the exploration-exploitation dilemma—that is, whether to (i) act greedily with the available data or to (ii) explore, which means taking sub-optimal—possibly unsafe—actions  $\mathbf{u}$  to learn a more accurate  $\hat{\mathbf{f}}$ .

RL typically assumes that the underlying control problem is a Markov decision process (MDP). An MDP comprises a state space  $\mathbb{X}$ , an input (*action*) space  $\mathbb{U}$ , stochastic dynamics (also called *transition model*, see M1), and a per-step reward function (see M2). When all the components of an MDP are known (in particular,  $\mathbf{f}$  and  $J$  in Sec. 2.1), then it solves the problem in Equations 9a, 9b, 9c without the constraints. Dynamic programming (DP) algorithms such as value and policy iteration can be used to find an optimal policy  $\pi^*$ . Many RL approaches, however, make no assumptions on any part of  $\mathbf{f}$  being known *a priori*.

We can distinguish model-based RL approaches, which learn an explicit model  $\hat{\mathbf{f}}$  of the system dynamics  $\mathbf{f}$  and use it to optimize a policy, from model-free RL algorithms. The latter algorithms (29) can be broadly categorized as: (i) value function-based methods, learning a value function (see M8); (ii) policy-search and policy-gradient methods, directly trying to find an optimal policy  $\pi^*$ ; and (iii) actor-critic methods, learning both a value function (*critic*) and the policy (*actor*). We also note that the convergence of these methods has been shown for simple scenarios but is still a challenge for more complex scenarios or when function approximators are used (30).

There are multiple practical hurdles to the deployment of RL algorithms in real-world robotics problems (6), for example, (i) the continuous, possibly high-dimensional  $\mathbb{X}$  and  $\mathbb{U}$  in robotics (often assumed to be finite, discrete sets in RL), (ii) the stability and convergence of the learning algorithm (31) (necessary, albeit not sufficient, to produce a stable policy, see M3), (iii) learning robust policies from limited samples, (iv) the interpretability of the learned policy—especially in deep RL when leveraging neural networks (NNs) (see M9) for function approximation (29)—and, importantly, (v) providing provable safety guarantees.

The exploration-exploitation dilemma can be mitigated using Bayesian inference. This



is achieved by computing posterior distributions over the states  $\mathbb{X}$ , possible dynamics  $\mathbf{f}$ , or total cost  $J$  from past observed data (32). These posteriors provide explicit knowledge of the problem’s uncertainty and can be used to guide exploration. In practice, however, full posterior inference is almost always prohibitively expensive and concrete implementations must rely on approximations.

To achieve constraint satisfaction (over states or sequences of states) and robustness to different, possibly noisy dynamics  $\mathbf{f}$ , constrained MDPs and robust MDPs are extensions of traditional MDPs that more closely resemble the problem statement in Sec. 2.1.

**Constrained MDPs** (CMDPs) (33) extend simple MDPs with constraints and optimize the problem in Equation 9 when Equation 9e takes the discounted form of *Safety Level  $I$* ’s Equation 6. We refer to the discounted constraint cost  $J_{c^j}$  under policy  $\pi$  as  $J_{c^j}^\pi$ . Traditional approaches to solve CMDPs, such as Linear Programming and Lagrangian methods, often assume discrete state-action spaces and cannot properly scale to complex tasks such as robot control. Deep RL promises to mitigate this, yet, applying it to the constrained problem still suffers from the computational complexity of the off-policy evaluation (M10) of trajectory-level constraints  $J_{c^j}$  (34). In Sec. 3.2.3, we present some recent advances in CMDP-based work that feature: (i) integration of deep learning techniques in CMDPs for more complex control tasks; (ii) provable constraint satisfaction throughout the exploration or learning process; and (iii) constraint transformation for the efficient evaluation of  $J_{c^j}$  from data collected off-policy.

**Robust MDPs** (35), inspired by robust control (see Sec. 2.2), extend MDPs such that the dynamics can include parametric uncertainties or disturbances, and the cost of the worst-case scenario is optimized. This is captured by the min-max optimization problem:

$$J^\pi(\bar{\mathbf{x}}_0) = \min_{\pi_{0:N-1}} \max_{\hat{\mathbf{f}} \in \mathbb{D}} J(\mathbf{x}_{0:N}, \mathbf{u}_{0:N-1}) \quad 13a.$$

$$\text{s.t. Equation 9b, Equation 9c, Equation 9d,} \quad 13b.$$

where  $\mathbb{D}$  is a given uncertainty set of  $\hat{\mathbf{f}}$ . To keep solutions tractable, practical implementations typically restrict  $\mathbb{D}$  to certain classes of models. This can limit the applicability of robust MDPs beyond toy problems. Recent work (36, 37, 38) applied deep RL to robust decision making, targeting key theoretical and practical hurdles such as (i) how to effectively model uncertainty with deep neural networks, and (ii) how to efficiently solve the min-max optimization (e.g., via sampling or two-player, game-theoretic formulations). These ideas, including adversarial RL and domain randomization, are presented in Sec. 3.2.4.

## 2.4. Bridging Control Theory and RL for Safe Learning Control

When designing a learning-based controller, we typically have two sources of information: (i) our prior knowledge, and (ii) data generated by the robot system. Control approaches rely on prior knowledge and on assumptions such as parametric dynamics models to provide safety guarantees. RL approaches typically make use of expressive learning models to extract patterns from data that facilitate learning complex tasks, but these can impede the provision of any formal guarantees. In recent literature, we see an effort from both the control and RL communities to develop safe learning control algorithms, with the goal of systematically leveraging expressive models for closed-loop control (see Figure 4). Questions that arise from these efforts include: how can control-theoretic tools be applied to expressive machine learning models? Or, how can expressive models be incorporated into control frameworks?

Expressive learning models can be categorized as deterministic (e.g., standard deep neural networks (DNNs)) or probabilistic (e.g., Gaussian processes (GPs), see M11, and Bayesian linear regression (BLR)). Deep learning techniques such as feed-forward NNs, convolutional NNs, and long short-term memory networks have the advantage of being able

---

### M10. Off-policy, On-policy, and Offline RL:

Off-policy RL improves the value functions estimates with data collected operating under different policies. On-policy methods only use data generated by the current policy. Offline RL uses data previously collected by an unknown policy.

---



---

### M11. Gaussian Process:

A *Gaussian process* (GP) is a set of random variables in which any finite number of the random variables can be modeled as a joint Gaussian distribution. Practically, a GP is a probabilistic model specifying a distribution over functions.

---

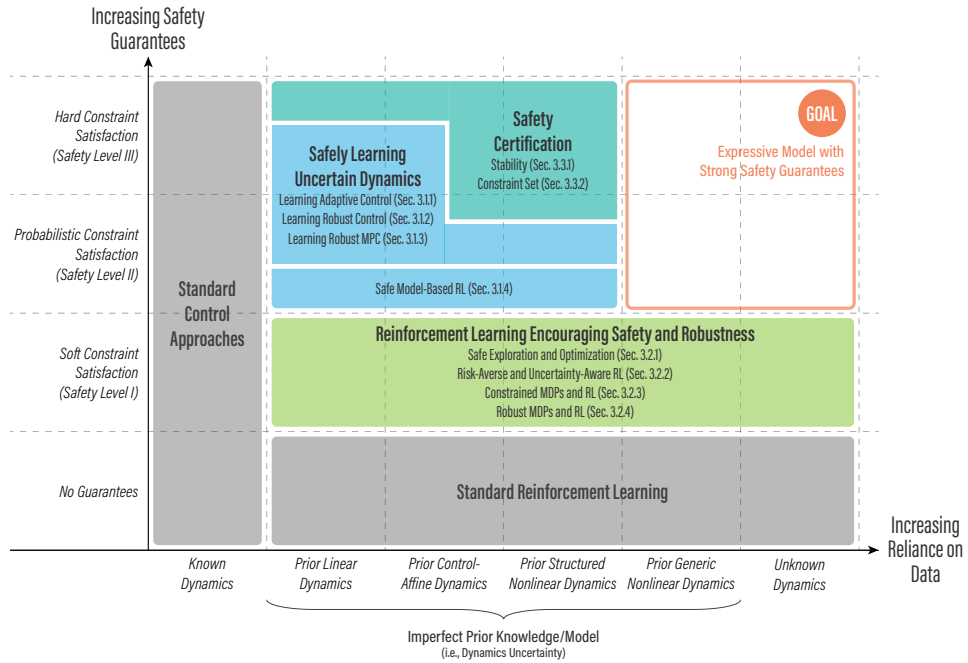


Figure 4: A summary of the safe learning control approaches reviewed in [Sec. 3](#).

to abstract large volumes of data, enabling real-time execution in a control loop. On the other hand, their probabilistic counterparts, such as GPs and BLRs, provide model output uncertainty estimates that can be naturally blended into traditional adaptive and robust control frameworks. We note that there are approaches aiming to combine the advantages of the two types of learning (e.g., Bayesian NNs), and quantifying uncertainty in deep learning is still an active research direction.

In this review, we focus on approaches that address the problem of safe learning control at two stages: (i) *online adaptation or learning*, where online data is used to adjust the parameters of the controller, the robot dynamics model, the cost function, or the constraint functions during closed-loop operation; and (ii) *offline learning*, where data collected from each trial is recorded and used to update a model in a batch manner in between trials of closed-loop operation. In safe learning control, data is generally used to address the issue of uncertainties in the problem formulation and reduce the conservatism in the system design, while the safety aspect boils down to “knowing what is not known” and cautiously accounting for the incomplete knowledge via algorithm design.

### 3. SAFE LEARNING CONTROL APPROACHES

The ability to guarantee safe robot control is inevitably dependent on the amount of prior knowledge available and the type of uncertainties present in the problem of interest. In this section, we discuss approaches for *safe learning control in robotics* based on the following three categories (see [Figure 2](#) and [Figure 4](#)):

- **Learning uncertain dynamics to safely improve performance:** The first set of works we review relies on an *a priori* model of the robot dynamics. The robot’s performance is improved by learning the uncertain dynamics from data. Safety is typically guaranteed based on standard control-theoretic frameworks, achieving Safety

Level II or III.

- **Encouraging safety and robustness in RL:** The second set of works encompasses approaches that usually do not have knowledge of an *a priori* robot model or the safety constraints. Rather than providing hard safety guarantees, these approaches encourage safe robot operation (Safety Level I), for example, by penalizing dangerous actions.
- **Certifying learning-based control under dynamics uncertainty:** The last set of works aims to provide safety certificates for learning-based controllers that do not inherently consider safety constraints. These approaches modify the learning controller output by either constraining the control policy, leveraging a known safe backup controller, or modifying the controller output directly to achieve stability and/or constraint satisfaction. They typically achieve Safety Level II or III.

**Figure 4** categorizes the approaches reviewed in this section based on their safety level and reliance on data. A more detailed summary of the approaches can be found in **Table 1** of **Appendix A**.

### 3.1. Learning Uncertain Dynamics to Safely Improve Performance

In this section, we consider approaches that improve the robot’s performance by learning the uncertain dynamics and provide safety guarantees (Safety Level II or III) via control frameworks such as adaptive control, robust control, and robust MPC (outlined in **Sec. 2.2**). These approaches make assumptions about the unknown parts of the problem (e.g., Lipschitz continuity) and often rely on a known model structure (e.g., control-affine or linear system with bounded uncertainty) to prove stability and/or constraint satisfaction (see **Figure 4**).

**3.1.1. Integrating Machine Learning and Adaptive Control.** There are three main ideas to incorporating online machine learning into traditional adaptive control (see **Sec. 2.2**), each with its own distinct benefits. These include: (i) using black-box machine learning models to accommodate nonparametric unknown dynamics, (ii) using probabilistic learning and explicitly accounting for the learned model uncertainties, in order to achieve cautious adaptation, and (iii) augmenting adaptive control with deep learning approaches for experience memorization in order to minimize the need for re-adaptation.

**Learning Nonparametric Unknown Dynamics with Machine Learning Models.** One goal of integrating adaptive control and machine learning is to improve the performance of a robot subject to nonparametric dynamics uncertainties. As opposed to **Equation 10**, we consider a system with nominal dynamics  $\bar{\mathbf{f}}(\mathbf{x}_k, \mathbf{u}_k) = \bar{\mathbf{A}}\mathbf{x}_k + \bar{\mathbf{B}}\mathbf{u}_k$  and unknown dynamics  $\hat{\mathbf{f}}(\mathbf{x}_k, \mathbf{u}_k) = \bar{\mathbf{B}}\psi_k(\mathbf{x}_k)$ , where  $\psi_k(\mathbf{x}_k)$  is an unknown nonlinear function without an obvious parametric structure. Learning-based MRAC approaches (**39, 40**) make the uncertain system behave as the linear nominal model  $\bar{\mathbf{f}}$  by using a combination of  $\mathcal{L}_1$  adaptation (**41**) and online learning to approximate  $\psi_k(\mathbf{x}_k)$  by  $\hat{\psi}_k(\mathbf{x}_k) = \hat{\psi}_{\mathcal{L}_1,k} + \hat{\psi}_{\text{learn},k}(\mathbf{x}_k)$ , where  $\hat{\psi}_{\mathcal{L}_1,k}$  is the input computed by the  $\mathcal{L}_1$  adaptive controller and  $\hat{\psi}_{\text{learn},k}$  is the input computed by the learning module, which can be an NN (**39**) or a GP (**40**). The estimated  $\hat{\psi}_k(\mathbf{x}_k)$  is then used in the controller  $\boldsymbol{\pi}_k(\mathbf{x}_k)$  to account for the unknown nonlinear dynamics  $\psi_k(\mathbf{x}_k)$ , improving a linear nominal control policy designed based on  $\bar{\mathbf{f}}$ . In these approaches, fast adaptation and stability guarantees are provided by the  $\mathcal{L}_1$  adaptation framework (Safety Level III), while the learning module provides additional flexibility to capture the unknown dynamics. As shown in (**39**), the addition of learning improves the performance of the standard  $\mathcal{L}_1$  adaptive controller and allows fast adaptation to be achieved at a lower sampling rate.

**Cautious Adaptation with Probabilistic Model Learning.** Another stream of adap-

tive control approaches leverage probabilistic models to achieve cautious adaptation by weighting the contribution of the learned model based on the model output uncertainty. We consider a system with nominal dynamics  $\bar{\mathbf{f}}(\mathbf{x}_k, \mathbf{u}_k) = \bar{\mathbf{f}}_{\mathbf{x}}(\mathbf{x}_k) + \bar{\mathbf{f}}_{\mathbf{u}}(\mathbf{x}_k)\mathbf{u}_k$  and unknown dynamics of the same form  $\hat{\mathbf{f}}(\mathbf{x}_k, \mathbf{u}_k) = \hat{\mathbf{f}}_{\mathbf{x}}(\mathbf{x}_k) + \hat{\mathbf{f}}_{\mathbf{u}}(\mathbf{x}_k)\mathbf{u}_k$ , where  $\hat{\mathbf{f}}_{\mathbf{x}}(\mathbf{x}_k)$  and  $\hat{\mathbf{f}}_{\mathbf{u}}(\mathbf{x}_k)$  are unknown nonparametric nonlinear functions. In a model inversion-based MRAC framework, an approximate feedback linearization is achieved via the nominal model to facilitate the design of MRAC, and a GP-based adaptation approach is used to compensate for feedback linearization errors due to the unknown dynamics (42). To account for the uncertainty in the GP model learning, the controller relies on the GP model only if the confidence in the latter is high:  $\boldsymbol{\pi}_k(\mathbf{x}_k) = \boldsymbol{\pi}_{\text{nom}}(\mathbf{x}_k) - \gamma(\mathbf{x}_k, \mathbf{u}_k) \boldsymbol{\pi}_{\text{learn},k}(\mathbf{x}_k)$ , where  $\boldsymbol{\pi}_{\text{nom}}(\mathbf{x}_k)$  is the control policy designed based the nominal model,  $\boldsymbol{\pi}_{\text{learn},k}(\mathbf{x}_k)$  is the adaptive component designed based on the GP model, and  $\gamma(\mathbf{x}_k, \mathbf{u}_k) \in [0, 1]$  is a scaling factor with 0 indicating low confidence in the GP. The stability of the overall system (Safety Level III) is guaranteed via a stochastic stability analysis, and the efficacy of the approach has been demonstrated in quadrotor experiments (43, 42).

**Memorizing Experience with Deep Architectures.** Apart from compensating for nonlinear and nonparametric dynamics uncertainties, deep learning approaches have also been applied to adaptive control for “memorizing” generalizable feature functions as the system adapts. In particular, an asynchronous DNN adaptation approach is proposed in (44, 45). Similar to (39, 40), we consider a linear nominal model  $\bar{\mathbf{f}}(\mathbf{x}_k, \mathbf{u}_k) = \bar{\mathbf{A}}\mathbf{x}_k + \bar{\mathbf{B}}\mathbf{u}_k$  and unknown nonlinear dynamics  $\hat{\mathbf{f}}(\mathbf{x}_k, \mathbf{u}_k) = \bar{\mathbf{B}}\psi_k(\mathbf{x}_k)$  with  $\psi_k(\mathbf{x}_k)$  being an unknown nonlinear function. In the proposed approach, the last layer of the DNN is updated at a higher frequency for fast adaptation, while the inner layers are updated at a lower frequency to “memorize” pertinent features for the particular operation regimes. To provide safety guarantees, the authors derive an upper bound on the sampling complexity of the DNN to achieve a prescribed level of modelling error and leverage this result to show that Lyapunov stability of the adapted system can be guaranteed (Safety Level III) by ensuring the modelling error of the DNN is lower than a given bound. In contrast to other MRAC approaches, which do not usually retain a memory of the past experience, the inner layers of the asynchronous DNN in (44, 45) store relevant features that facilitate adaptation when similar scenarios arise in the future.

**3.1.2. Learning-based Robust Control.** Learning-based robust control improves the performance of classical robust control in Sec. 2.2 by using data to improve the linear dynamics model and reduce the uncertainty in Equation 11.

**Using a Gaussian Process Dynamics Model for Linear Robust Control.** The conservative performance of robust control in Sec. 2.2 is improved by updating the linear dynamics model and uncertainty in Equation 11 with a Gaussian Process (46) (GP). The unknown nonlinear dynamics  $\hat{\mathbf{f}}(\mathbf{x}_k, \mathbf{u}_k)$  are learnt as a GP, which is then linearized about an operating point. By linearizing the GP (as opposed to directly fitting a linear model) data close to the operating point are prioritized. The uncertain linear dynamics in Equation 11 are assumed to be modelled as  $\hat{\mathbf{A}} = \mathbf{A}_0 + \tilde{\mathbf{A}}\boldsymbol{\Delta}$ ,  $\hat{\mathbf{B}} = \mathbf{B}_0 + \tilde{\mathbf{B}}\boldsymbol{\Delta}$ , where  $\mathbf{A}_0$  and  $\mathbf{B}_0$  are obtained from the linearized GP mean,  $\tilde{\mathbf{A}}$  and  $\tilde{\mathbf{B}}$  are obtained from the linearized GP variance (often two standard deviations), and  $\boldsymbol{\Delta}$  represents a matrix with elements taking any value in the range of  $[-1, +1]$ . Further performance improvement is achieved by modelling  $\tilde{\mathbf{A}}$  and  $\tilde{\mathbf{B}}$  as state dependent (47). Additionally, (48) was able to achieve better performance than (46) by leveraging the GP’s distribution, while maintaining the same level of safety. The main advantage of these approaches is that they can robustly guarantee Safety Level II stability while improving performance. This is achieved by shrinking the GP uncertainty as more data is added, which improves the linear model  $\mathbf{A}_0$  and  $\mathbf{B}_0$  and reduces the uncertain component  $\tilde{\mathbf{A}}$  and  $\tilde{\mathbf{B}}$ . This has been shown on a quadrotor (46). However, these approaches

are limited to stabilization tasks and do not account for state and input constraints.

**Exploiting Feedback Linearization for Robust Learning-Based Tracking.** Trajectory tracking convergence, as opposed to the simpler stabilization task performed in (46), is guaranteed by exploiting the special structure of exactly feedback linearizable systems (49). This structure assumes that the nonlinear system dynamics in Equation 1 can be described by a linear nominal model, where  $\bar{\mathbf{A}}$  and  $\bar{\mathbf{B}}$  have an integrator chain structure. It also assumes that the unknown dynamics are  $\hat{\mathbf{f}}(\mathbf{x}_k, \mathbf{u}_k) = \bar{\mathbf{B}}\psi(\mathbf{x}_k, \mathbf{u}_k)$ , where  $\psi(\mathbf{x}_k, \mathbf{u}_k)$  is an unknown invertible function. A probabilistic upper bound is obtained for  $\psi(\mathbf{x}_k, \mathbf{u}_k)$  by learning this function as a GP. A robust linear controller is designed for the uncertain system based on this learnt probabilistic bound. Further performance improvement is achieved by also updating the feedback linearization (50) through improving the estimate of the inverse  $\psi^{-1}(\cdot)$ . These approaches have been applied to trajectory tracking, with Safety Level II, on Lagrangian mobile manipulators in (49) and quadrotor models in (50). However, they hinge on this special structure and cannot account for state and input constraints.

**3.1.3. Reducing Conservatism in Robust MPC with Learning and Adaptation.** The conservative nature of robust MPC (Sec. 2.2) is improved—while still satisfying input and state constraints—through (i) robust adaptive MPC, which adapts to parametric uncertainty, and (ii) learning-based robust MPC, which learns the unknown dynamics  $\hat{\mathbf{f}}$  or, in one case, cost  $\hat{J}$ .

**Robust Adaptive MPC.** Robust adaptive MPC assumes parametric uncertainties, and (i) uses data to reduce the set of possible parameters over time or (ii) uses an inner-loop adaptive controller and an outer-loop robust MPC. This leads to improved performance compared to standard robust MPC (see Sec. 2.2), while satisfying hard constraints (Safety Level III).

The first set of approaches considers stabilization tasks, where the full system dynamics, Equation 10, are assumed to be linear and stable (51) or linear and unstable (52) with uncertain parameter  $\theta \in \Theta_0$ :

$$\mathbf{x}_{k+1} = (\bar{\mathbf{A}} + \hat{\mathbf{A}}(\theta))\mathbf{x}_k + (\bar{\mathbf{B}} + \hat{\mathbf{B}}(\theta))\mathbf{u}_k + \mathbf{w}_k. \quad 14.$$

The process noise and the parameters are assumed to be bounded by known sets  $\mathbb{W}$  and the initially conservative  $\Theta_0$ , respectively. Given  $\mathbb{W}$  and  $\Theta_0$ , we can derive a conservative upper bound on the uncertain dynamics  $\hat{\mathbf{f}}(\mathbf{x}_k, \mathbf{u}_k, \mathbf{w}_k, \theta) = \hat{\mathbf{A}}(\theta)\mathbf{x}_k + \hat{\mathbf{B}}(\theta)\mathbf{u}_k + \mathbf{w}_k \in \mathbb{D}_0$ , where  $\mathbb{D}_0$  is a compact set determined from  $\mathbb{W}$  and  $\Theta_0$ . To guarantee constraint satisfaction, tube-based MPC (see Sec. 2.2) is applied, where the initial tube  $\Omega_{\text{tube},0}$  is based on  $\mathbb{D}_0$ . To reduce the conservatism of the approach, an adaptive control method is introduced to improve the estimate of the parameter set  $\Theta_k$  and reduce the size of the tube  $\Omega_{\text{tube},k}$  at each time step  $k$ . This idea has been extended to stochastic process noise for probabilistic constraint satisfaction (Safety Level II) (53), time-varying parameters (54), and linearly parameterized uncertain nonlinear systems (55, 56). Further performance improvement is achieved by combining robust adaptive MPC with iterative learning (57), which updates the terminal constraint set  $\mathbb{X}_{\text{term}}$  in Equation 12d and the terminal cost  $l_H$  in Equation 12a after every full iteration using the closed-loop state trajectory and cost (58). In the second approach, an underlying MRAC (see Sec. 2.2) is used to make the closed-loop system dynamics resemble a linear reference model with bounded disturbance set  $\mathbb{D}$  (59). This linear model and its bounds are then used in an outer-loop robust MPC to achieve fast stabilization in the presence of model errors.

**Learning-Based Robust MPC.** Learning-based robust MPC uses data to (i) improve the unknown dynamics estimate, (ii) reduce the uncertainty set, or (iii) update the cost to avoid states with high uncertainty. Unlike robust adaptive control, learning-based robust MPC considers nonparameterized systems.

Under the assumption of a linear nominal model  $\bar{\mathbf{f}}(\mathbf{x}, \mathbf{u}) = \bar{\mathbf{A}}\mathbf{x} + \bar{\mathbf{B}}\mathbf{u}$  and bounded unknown dynamics  $\hat{\mathbf{f}}(\mathbf{x}_k, \mathbf{u}_k) \in \mathbb{D}$ , the unknown dynamics can be safely learned from data, e.g., using a NN (60). Robust constraint satisfaction, Equation 12c, is guaranteed by using tube-based MPC for the linear nominal model and performance improvement is achieved by optimizing over the control inputs for the combined nominal and learned dynamics. If the bounded unknown dynamics are assumed to be state-dependent  $\hat{\mathbf{f}}(\mathbf{x}_k) \in \mathbb{D}(\mathbf{x}_k)$ , instead of using a constant tube  $\Omega_{\text{tube}}$ , the state constraints in Equation 12c can be tightened based on the state-dependent uncertainty set  $\mathbb{D}(\mathbf{x})$  (61). In a numerical stabilization task, a GP is used to model  $\hat{\mathbf{f}}(\mathbf{x}_k)$  and its covariance determines  $\mathbb{D}(\mathbf{x})$ . This yields less conservative, probabilistic state constraints (Safety Level II).

Learning-based robust MPC can be extended to nonlinear nominal models. Typically, a GP is used to learn the unknown dynamics  $\hat{\mathbf{f}}(\mathbf{x}, \mathbf{u})$ : the mean updates the dynamics model in Equation 12b and the state- and input-dependent uncertainty set  $\mathbb{D}(\mathbf{x}, \mathbf{u})$  is derived from the GP's covariance and contains the true uncertainty with high probability. Similarly, the state- and input-dependent tube in Equation 12c is determined from the uncertainty set  $\mathbb{D}(\mathbf{x}, \mathbf{u})$ . The main challenge, compared to using a linear nominal model, is the uncertainty propagation over the prediction horizon in the MPC because Gaussian uncertainty (obtained from the GP) is no longer Gaussian when propagated through the nominal nonlinear dynamics. Approximation schemes are required, such as using a sigma-point transform (62), linearization (63), exact moment matching (64), or ellipsoidal uncertainty set propagation (65). Additionally, further approximations (e.g., fixing the GP's covariances over the prediction horizon (63)) are usually required to achieve real-time implementation. These approximations can lead to violations of the probabilistic constraints (Safety Level II) in Equation 9e. An alternative approach to address the challenges of GPs uses a NN regression model that predicts the quantile bounds of the tail of a state trajectory distribution for a tube-based MPC (66). However, this approach is hindered by typically nonexhaustive training data sets, which can lead to the underestimation of the tube size.

For repetitive tasks, another approach is to adjust the cost function based on data instead of the system model. The predicted cost error is learned from data using the difference between the predicted cost at each time step and the actual closed-loop cost at execution. By adding this additional term to the cost function, the MPC penalizes states that had previously resulted in higher closed-loop cost than expected Equation 12a (67), resulting in reliable performance despite model errors.

**3.1.4. Safe Model-Based RL with A Priori Dynamics.** Safe model-based RL augments model-based RL (see Sec. 2.3) with safety guarantees (see Figure 4). Stability can be probabilistically guaranteed (Safety Level II) under the assumption that the known nominal model  $\bar{\mathbf{f}}(\mathbf{x}_k, \mathbf{u}_k)$  and the unknown part  $\hat{\mathbf{f}}(\mathbf{x}_k, \mathbf{u}_k)$  are Lipschitz continuous with known Lipschitz constants (see M12) (68). Given a Lyapunov function (see M6), an initial safe policy  $\pi_0$ , and a GP to learn the unknown dynamics  $\hat{\mathbf{f}}$ , a control policy  $\pi_k$  is chosen so that it maximizes a conservative estimate of the *region of attraction (ROA)* (see M13). The most uncertain states (based on the GP's covariance) inside the ROA, are explored. This reduces the uncertainty over time and allows the ROA to be extended. The practical implementation resorts to discrete states for tractability and retains the stability guarantees while being sub-optimal in (exploration) performance.

## 3.2. Encouraging Safety and Robustness in Reinforcement Learning

The approaches in this section are safety-augmented variations of the traditional MDP and RL frameworks. In general, these methods do not assume knowledge of an *a priori* nominal model  $\bar{\mathbf{f}}$  and some also learn the reward or step cost  $l$  (69), or the safety constraints  $\mathbf{c}$  (70). Rather than providing strict safety guarantees, these approaches encourage constraint satis-

### M12. Lipschitz

**Continuity:** A function  $\mathbf{h} : \mathbb{A} \mapsto \mathbb{B}$  is said to be *Lipschitz continuous* on  $\mathbb{A}$  if  $(\exists \rho > 0)$  such that for any  $\mathbf{a}_1, \mathbf{a}_2 \in \mathbb{A}$  the condition  $\|\mathbf{h}(\mathbf{a}_1) - \mathbf{h}(\mathbf{a}_2)\| \leq \rho \|\mathbf{a}_1 - \mathbf{a}_2\|$  holds, where  $\mathbb{A}$  and  $\mathbb{B}$  are the domain and codomain of  $\mathbf{h}$ , and  $\|\cdot\|$  denotes any  $l_p$ -norm of a vector. The constant  $\rho$  is called a *Lipschitz constant* of  $\mathbf{h}$ .

### M13. Region of

**Attraction:** Let  $\mathbf{x}^*$  be an equilibrium of a closed-loop system under some policy  $\pi(\mathbf{x})$ :  $\mathbf{x}_{k+1} = \mathbf{f}_\pi(\mathbf{x}_k) = \mathbf{f}(\mathbf{x}_k, \pi(\mathbf{x}_k))$  with  $\mathbf{f}_\pi(\mathbf{x}_k)$  being Lipschitz continuous. The *region of attraction (ROA)* of  $\mathbf{x}^*$  is the set of states from which the system will converge to  $\mathbf{x}^*$ .

faction during and after learning, or robustness of the learned control policy  $\pi$  to uncertain dynamics (Safety Level I, see **Figure 4**). In plain MDP formulations, (i) states and inputs (or *actions*) are assumed to have known, often discrete and finite, domains but they are not further constrained while searching for an optimal policy  $\pi^*$  and (ii) only loose assumptions are made on the dynamics  $\mathbf{f}$ , for example, the system being Markovian (20) (see M14).

A previous taxonomy of safe RL (4)—covering research published up until 2015—distinguished methods that either modify the exploration process to include external knowledge or modify the optimality criterion  $J$  with a safety factor. However, as pointed out in (71), the number and breadth of publications in RL, including safe RL, has since greatly increased. Because the approaches recently proposed in this area are numerous and diverse, we provide a high-level review of some of the most significant trends with a focus on their applicability to robotics. These include: (i) the safe exploration of MDPs; (ii) risk-aware and cautious RL adaptation; (iii) RL based on constrained MDPs; and (iv) robust RL.

**3.2.1. Safe Exploration and Optimization.** An RL algorithm’s need to explore poses a challenge to its safety as it must select inputs with unpredictable consequences in order to learn about them.

**Safe Exploration.** The problem of safely exploring an MDP is tackled in (73) through the notion of ergodicity (see M15). Policy updates that preserve the ergodicity of the MDP enable the system to return to an arbitrary state from any state. Thus, the core idea is to restrict the space of eligible policies to those that make the MDP ergodic (with at least a given probability). Exactly solving this problem, however, is NP-hard. A simplified problem is solved in (73), using a heuristic exploration algorithm (74), which leads to sub-optimal but safe exploration that only considers a subset of the ergodic policies. In practice, (73) is demonstrated in two simulated scenarios with discrete, finite  $\mathbb{X}$  and  $\mathbb{U}$ . Rather than designing for a recoverable exploration through ergodicity, safe exploration strategies that provide constraint satisfaction are developed in (75, 70). A *safety layer* is used to convert an optimal but potentially unsafe action  $\mathbf{u}_{\text{learn},k}$  produced by an NN policy, to the closest safe action  $\mathbf{u}_{\text{safe},k}$  with respect to some safety state constraint. Both (75) and (70) involve solving a constrained least squares problem,  $\mathbf{u}_{\text{safe},k} = \arg \min_{\mathbf{u}_k} \|\mathbf{u}_k - \mathbf{u}_{\text{learn},k}\|_2^2$ , which is akin to the safety filter approaches further discussed in Sec. 3.3.2. Concretely, (75) assumes full knowledge of the (linear) constraint functions and solves  $\mathbf{u}_{\text{safe},k}$  using a differentiable Quadratic Programming (QP) solver. In contrast, (70) assumes constraints are unknown *a priori* but can be evaluated. Thus, the approach learns the parameters of linear approximations to these constraints and then uses them in the solver. Notably, (75, 70) consider single time-step state constraints. In Sec. 3.2.3, we will also discuss methods that deal with more general trajectory-level constraints (Equation 6).

**Safe Optimization.** Several works address the problem of safely optimizing an unknown function (typically the cost function), often exploiting GP models (76). Safety refers to sampling inputs that do not violate a given safety threshold (Equation 3, Equation 4). These approaches fall under the category of Bayesian optimization (77) and include *SafeOpt* (78); *SafeOpt-MC* (79), an extension towards multiple constraints; *StageOpt* (80), a the more efficient two-stage implementation; and *GoSafe* (81), to explore beyond the initial safe region. In particular, *SafeOpt* infers two subsets of the safe set from the GP model: (i) one with candidate inputs to extend the safe set, and (ii) one with candidate input to optimize the unknown function—from which it greedily picks the most uncertain. The ideas pioneered by *SafeOpt* were applied to MDPs in *SafeMDP* (69) resulting in safe exploration of MDPs with an unknown cost function  $l(\mathbf{x}, \mathbf{u})$ , which (69) models as a GP. In *SafeMDP*, the single-step reward represents the safety feature that should not violate a given threshold. Another extension of (78), *SafeExpOpt-MDP* (82), treats the safety feature  $c^j$  and the MDP’s cost  $l$  as two separate, unknown functions, allowing for the constraint of the first and the optimization

---

**M14. Markov Property:** The probability for a system to be in state  $\mathbf{x}_k$  at time  $k$  only depends on the state at time  $k - 1$ ,  $\mathbf{x}_{k-1}$  and, in MDPs,  $\mathbf{u}_{k-1}$ .

---

**M15. Ergodicity:** An MDP is *ergodic* if, by following a given policy, any state is reachable from any other state. Ergodic properties of dynamical systems are key to establish conditions for their controllability (72).

---

---

**M16. Conservative Q-learning (CQL):**

A recent offline RL approach (83) to mitigate the overestimation problem in value functions. Q-learning is known to overestimate action-state values because of its maximization operator, and offline RL can lead to overestimated values of the actions that are more likely under the policy used to collect the data. CQL tackles this by using a novel update rule with a simple regularizer to learn value lower bounds instead.

---

**M17. Model Ensembles:**

Collections of learned models can be used to mitigate noise or better capture the stochasticity of a system. *Probabilistic Ensembles with Trajectory Sampling* (PETS) (84) is a model-based RL method that uses model ensembles and probabilistic networks to capture uncertainties in the system dynamics.

---

**M18. Conditional Value at Risk (CVaR):**

$\text{CVaR}_\alpha$  is a risk measure that is equal to the average of the worst  $\alpha$ -percentile of the total cost  $J$  (88).

---

of the latter. A recent survey of these techniques was provided in (76), highlighting the distinction between safe learning in regression—i.e., minimizing the selection and evaluation of non-safe training inputs—and safe exploration in MDPs and stochastic dynamical systems like Equation 9b—i.e., selecting action inputs that also preserve ergodicity).

**Learning a Safety Critic.** A safety critic is a learnable action-value function  $Q_{\text{safe}}^\pi$  that can be used to detect if a proposed action can lead to unsafe conditions. The works (85, 86, 87) make use of such a critic and then resort to various fallback schemes to determine a safer alternative input. These works differ from those in Sec. 3.3.2 in that the filtering criterion depends on a model-free, learned value function, which can only grant the satisfaction of Safety Level I.

In (85), *safety Q-functions for reinforcement learning* (SQRL) are used to (i) learn a safety critic from only abstract, sparse safety labels (e.g., a binary indicator), and (ii) transfer knowledge of safe action inputs to new but similar tasks. SQRL trains  $Q_{\text{safe}}^\pi$  to predict the future probability of failure in a trajectory, and uses it to filter out unsafe actions from the policy  $\pi$ . The knowledge transfer is achieved by pre-training  $Q_{\text{safe}}^\pi$  and  $\pi$  in simulations, and then fine-tuning  $\pi$  on the new task (with similar dynamics  $\mathbf{f}$  and safety constraints), still in simulation, while reusing  $Q_{\text{safe}}^\pi$  to discriminate unsafe inputs. However, the success of the final safe policy still depends on the task- and environment-specific hyperparameters (85) (which must be found *via* parameter search, in simulation, prior to the actual experiment). Building on (85), *Recovery RL* (86) additionally learns a recovery policy  $\pi_{\text{rec}}$  to produce fallback actions for  $Q_{\text{safe}}^\pi$  as an alternative to filtering out unsafe inputs and resorting to potentially sub-optimal ones in  $\pi$ . In (87), the authors extend Conservative Q-Learning (CQL, see M16) and propose Conservative Safety Critic (CSC). Similarly to (85), CSC assumes sparse safety labels and uses  $Q_{\text{safe}}^\pi$  for action filtering, but it trains  $Q_{\text{safe}}^\pi$  to upper bound the probability of failure and ensures provably safe policy improvement at each iteration.

**3.2.2. Risk-Averse RL and Uncertainty-Aware RL.** Safety in RL can also be encouraged by deriving and using risk (or uncertainty) estimates during learning. These estimates are typically computed for the system dynamics or the overall cost function, and leveraged to produce more conservative (and safer) policies.

Risk can be defined as the probability of collision for a robot performing a navigation task (89, 90). A collision model, captured by a neural network ensemble trained with Monte Carlo dropout, predicts the probability distribution of a collision, given the current state and a sequence of future actions. The collision-averse behavior is then achieved by incorporating the collision model in a MPC planner.

In (91), the authors build upon PETS (84) (see M17) to propose *cautious adaptation for safe RL* (CARL). CARL is composed of two training steps: (i) a pre-training phase that is not risk-aware, where a PETS agent is trained on multiple different system dynamics, and (ii) an adaptation phase, where the agent is fine-tuned on the target system by taking risk-averse actions. CARL also proposed two notions of risk, one to avoid low-reward trajectories and another to avoid *catastrophes* (e.g., irrecoverable states or constraint violations). In SAVED (92), a PETS (84) agent is used to predict the probability of a robot's collision and to evaluate a chance constraint for safe exploration. Similarly to the methods in the last paragraph of Sec. 3.2.1, SAVED also learns a value function from sparse costs, which it uses as a terminal cost estimate.

To learn risk-averse policies using only offline data, the approach in (93) optimizes for a risk measure of the cost such as Conditional Value-at-Risk (CVaR) (see M18). Instead of using model ensembles as in (89, 90, 91), the work in (93) uses distributional RL (see M19) to explicitly model the distribution of the total cost of the task (control of a simulated 1D car), and offline learning to improve scalability.



**3.2.3. Constrained MDPs and RL.** The CMDP framework is frequently used in safe RL, as it introduces constraints that can be used to express arbitrary safety notions in the form of Equation 6. RL for CMDPs, however, faces two important challenges (see Sec. 2.3): (i) how to incorporate and enforce constraints in the RL algorithm; and (ii) how to efficiently solve the constrained RL problem—especially when using deep learning models, which are the *de facto* standard in non-safe RL. Works in this section that address these questions include (i) Lagrangian methods for RL, (ii) generalized Lyapunov functions for constraints, and (iii) Backward Value Functions. Nonetheless, much of the work in this section remains confined to rather naive simulated tasks, motivating further investigation into the applicability of constrained RL in real-world control.

**Lagrangian Methods in RL Optimization.** In (88, 95), the CMDP optimization problem (see Sec. 2.3) is first transformed into an unconstrained optimization one (see M20) over the primal variable  $\pi$  and dual variable  $\lambda$  using the Lagrangian  $\mathcal{L}(\pi, \lambda)$ , and RL is used as a subroutine in the primal-dual updates for Equation 15b:

$$\mathcal{L}(\pi, \lambda) = J^\pi + \sum_j \lambda_j \left( J_{c_j}^\pi - d^j \right) \quad 15a.$$

$$(\pi^*, \lambda^*) = \arg \max_{\lambda \geq 0} \min_{\pi} \mathcal{L}(\pi, \lambda) \quad \text{s.t.} \quad \text{Equation 9b, 9c, 9d.} \quad 15b.$$

In particular, (88) defines a constraint on the CVaR of cost  $J^\pi$ , and uses policy gradient or actor-critic methods to update the policy in Equation 15b. In subsequent work (95), the authors improve upon (88) by incorporating off-policy updates of the dual variable (with the on-policy primal-dual updates). This is empirically shown to achieve better sample efficiency and faster convergence. In (34), the authors extend a standard trust-region (see M21) RL algorithm (96) to CMDPs using a novel bound that relates the expected cost of two policies to their state-averaged divergence (see M21). The key idea is performing primal-dual updates with surrogates/approximations of the cost  $J^\pi$  and constraint cost  $J_{c_j}^\pi$  derived from the bound. The benefits are two-fold: (i) the surrogates can be estimated with only state-action data, bypassing the challenge of trajectory evaluation from off-policy data; and (ii) the updates guarantee monotonic policy improvement and near constraint satisfaction at each iteration (Safety Level I). However, unlike (88, 95), each update involves solving the dual variables from scratch, which can be computationally expensive.

**A Lyapunov Approach to Safe RL.** Lyapunov functions (see M6) are used extensively in control to analyze system stability and they are a powerful tool to translate a system’s global properties into local conditions. In (97, 98), Lyapunov functions are used to transform the trajectory-level constraints  $J_{c_j}^\pi$  in Equation 6 into step-wise, state-based constraints. This allows a more efficient computation of  $J_{c_j}^\pi$  and mitigates the cost of off-policy evaluation. These approaches, however, require the system to start from a baseline policy  $\pi_0$  that already satisfies the constraints. In (97), four different algorithms are proposed to solve CMDPs by combining traditional RL methods and the Lyapunov constraints, but they are only applicable to discrete input spaces (with continuous state spaces). In subsequent work (98), the authors of (97) extend the approach to continuous input spaces and standard policy gradient methods, addressing its computational tractability.

**Learning Backward Value Functions.** The work in (99) proposes Backward Value Functions (BVF) as a way to overcome the excessive computational cost of the approaches in the previous paragraphs (34, 98). Similarly to a (forward) value function  $V^\pi$  that estimates the total future cost from each state, a BVF  $V^{b,\pi}$  estimates the accumulated cost so far (up to the current state). We can decompose a trajectory-level constraint at any time step  $k$  as sum of  $V_{c_j}^\pi(\mathbf{x}_k)$  and  $V_{c_j}^{b,\pi}(\mathbf{x}_k)$  for the constraint cost  $J_{c_j}^\pi$ . This decomposition also alleviates the problem of off-policy evaluation as these value functions can be learned concurrently and efficiently via Temporal Difference (TD) methods (see M22). In practice,  $V_{c_j}^\pi, V_{c_j}^{b,\pi}$ ,

---

**M19. Distributional**

**RL:** Distributional RL (94) focuses on the modeling and learning of the *distribution* of the cost  $J$ , rather than its expected value. The greater descriptive power can enable richer predictions and risk-aware behaviours.

---

**M20. Lagrangian**

**Methods:** To solve a constrained optimization problem

$$\min_x f(x) \text{ s.t. } g(x) \leq 0, \text{ Lagrangian}$$

*methods* define a Lagrangian function

$$\mathcal{L}(x, \lambda) = f(x) + \lambda g(x) \text{ on}$$

primal variable  $x$ , dual variable  $\lambda$ , and solve the equivalent unconstrained problem

$$\max_{\lambda \geq 0} \min_x \mathcal{L}(x, \lambda).$$

In practice, numerical approaches, such as gradient methods, are often used to perform the primal-dual updates (88).

---

---

**M21. Divergence:**

*Divergence* quantifies the similarity between probability distributions. It may not satisfy the symmetry and the triangle inequality of a distance metric. In RL, commonly used ones include Kullback-Leibler divergence (KL), total variation distance (TV), and the Wasserstein distance, which is also a valid distance metric. Divergence measures can be used to constrain policy updates in trust-region RL algorithms (96).

---

**M22. Temporal Difference Learning:**

*Temporal Difference (TD) learning* (20) is an important class of model-free RL methods that updates the value function by bootstrapping from the current estimate using the Bellman equation to formulate sampling-based, iterative updates. For example, given a control step data tuple  $(\mathbf{x}_k, l_k, \mathbf{x}_{k+1})$  we can perform a learning step as  $V^\pi(\mathbf{x}_k) \leftarrow V^\pi(\mathbf{x}_k) + \alpha(l_k + \gamma V^\pi(\mathbf{x}_{k+1}) - V^\pi(\mathbf{x}_k))$  where  $\alpha, \gamma$  are the step size and discount factor.

---

$V^\pi$  are jointly learned (99) and used for policy improvement at each time step, allowing to implement Safety Level I. The approach is intended for discrete action spaces but can be adapted to continuous ones as in (70).

**3.2.4. Robust MDPs and RL.** Works in this section aim to implement robustness in RL, specifically, learning policies that can operate under disturbances and generalize across similar tasks or robotic systems. This is typically done by framing the learning problem as a robust MDP (Equation 13). In (100), Robust RL (RRL) implements an Actor-Disturber-Critic architecture. The authors observe that the learned policy and value function coincide with those derived analytically from  $H_\infty$  control theory for linear systems (see Sec. 2.2). However, more recent robust RL literature often abstains from assumptions on dynamics or disturbances, and applies model-free RL (101) to seek empirically robust performance at the expense of theoretical guarantees. In the following paragraphs, we introduce two lines of work: (i) robust adversarial RL, which explicitly models the min-max problem in Equation 13 in a game theoretic fashion; and (ii) domain randomization, which approximates the same problem in Equation 13 by learning on a set of randomized perturbed dynamics.

**Robustness Through Adversarial Training.** Combining RL with adversarial learning (102) results in robust adversarial RL (RARL) (36), where the robust optimization problem (Equation 13) is set up as a two-player, discounted zero-sum Markov game in which an agent (protagonist) learns policy  $\pi$  to control the system and another agent (adversary) learns a separate policy to destabilize the system. The two agents learn in an alternated fashion (each is updated while fixing the other), attempting to progressively improve both the robustness of the protagonist's policy and the strength of its adversary. The work of (37) extends (36) with risk-aware agents (see Sec. 3.2.2), with the protagonist being risk-averse and the adversary being risk-seeking. This method learns an ensemble of Deep Q-Networks (DQN) (103) and defines the risk of an action based on the variance of its value predictions. In another extension of (36), a population of adversaries (rather than a single one) is trained (38), leading to the resulting protagonist being less exploitable by new adversaries. Finally, the work in (104) proposes certified lower bounds for the value predictions from a DQN (103), given bounded observation perturbations. The action selection is based on these value lower bounds, assuming adversarial perturbation.

**Robustness Through Domain Randomization.** Domain randomization methods attempt to learn policies that empirically generalize to a wider range of tasks or robot systems. Instead of considering worst-case disturbances or scenarios, learning happens on systems with randomly perturbed parameters (e.g., inertial properties, friction coefficients). These parameters often have pre-specified ranges, effectively inducing a robust set that (i) approximates the uncertainty set  $\mathbb{D}$  in Equation 13, and (ii) allows the system to reuse any standard RL algorithm. As an example, in (105), a quadrotor learns vision-based flight purely in simulation from scenes with randomized properties. Using this model-free policy in the real world resulted in improved collision avoidance performance. Instead of learning a policy directly, the work in (106) uses learned visual predictions with a MPC controller to enable efficient and scalable real-world performance. Besides the uniform randomization in (105, 106), adaptive randomization strategies such as Bayesian search (107) are also a promising direction. In (108), a discriminator is adversarially trained and used to guide the randomization process that generates systems that are less explored or exploited by the current policy.

### 3.3. Certifying Learning-Based Control Under Dynamics Uncertainty

In this section, we review methods providing certification to learning-based control approaches that do not inherently account for safety constraints, see Figure 2. We divide the

discussion into two parts: (i) stability certification, and (ii) constraint set certification. The works in this section leverage an *a priori* dynamics model of the system and provide hard or probabilistic safety guarantees (Safety Level II or III) under dynamics uncertainties (see [Figure 4](#)).

**3.3.1. Stability Certification.** This section introduces certification approaches that guarantee closed-loop stability under a learning-based control policy.

***Lipschitz-Based Safety Certification for DNN-Based Learning Controllers.*** These approaches exploit the expressive power of DNNs for policy parametrization and guarantee closed-loop stability through a Lipschitz constraint on the DNN. Let  $\rho$  be a Lipschitz constant of a DNN policy (see [M12](#)), an upper bound on the Lipschitz constant  $\rho$  that guarantees closed-loop stability (Safety Level III) can be established by using a small-gain stability analysis ([109](#)), solving a semi-definite program (SDP) ([110](#)), or applying a sliding mode control framework ([111](#)). This bound on  $\rho$  can be used in either (i) a passive, iterative enforcement approach where the Lipschitz constant  $\rho$  is first estimated (e.g., via an SDP-based estimation ([112](#))) and then used to guide re-training until the Lipschitz constraint is satisfied, or (ii) an active enforcement approach where the Lipschitz constraint is directly enforced by the training algorithm of the DNN (e.g., via spectral normalization ([111](#))). While guaranteeing stability, the Lipschitz-based certification approaches often rely on the particular structure of the system dynamics (e.g., a control-affine structure or linear structure with additive nonlinear uncertainty) to find the certifying Lipschitz constant. It remains to be explored how this idea can be extended to more generic robot systems.

***Learning Regions of Attraction for Safety Certification.*** The region of attraction (ROA) of a closed-loop system (see [M13](#)) is used in the learning-based control literature as a means to guarantee safety. For a nonlinear system with a given state-feedback controller, the ROA is the set of states that is guaranteed to converge to the equilibrium, which is treated as a safe state. This notion of safety provides a mean to certify a learning-based controller. It guarantees that there is a region in state space from which the controller can drive the system back to the safe state (Safety Level III) ([113](#)). ROAs can, for example, be used to guide data acquisition for model or controller learning ([68, 114](#)).

We consider deterministic closed-loop systems  $\mathbf{x}_{k+1} = \mathbf{f}_\pi(\mathbf{x}_k) = \mathbf{f}(\mathbf{x}_k, \pi(\mathbf{x}_k))$  with  $\mathbf{f}_\pi(\mathbf{x}_k)$  being Lipschitz continuous. A Lyapunov neural network (LNN) can be used to iteratively learn the ROA of a controlled nonlinear system from the system’s input-output data ([113](#)). As compared to the typical Lyapunov functions in control (e.g., quadratic Lyapunov functions), the proposed method uses the LNN as a more flexible Lyapunov function representation to provide a less conservative estimate of the system’s ROA. The necessary properties of a Lyapunov function are preserved via the network’s architectural design. In ([114](#)), an ROA estimation approach for high-dimensional systems is presented, combining a sum of squares (SOS) programming method for the ROA computation ([115](#)) and a dynamics model order reduction technique to curtail computational complexity ([116](#)).

**3.3.2. Constraint Set Certification.** This section summarizes approaches that provide constraint set certification to a learning-based controller based on the notion of *robust positive control invariant* (RPCI) safe sets  $\Omega_{\text{safe}} \subseteq \mathbb{X}_c$  (see [M23](#)). Certified learning, which can be achieved through a safety filter ([7](#)) or shielding ([117](#)), finds the minimal modification of a learning-based control input  $\mathbf{u}_{\text{learn}}$  (see [Figure 2](#)) such that the system’s state stays inside

---

**M23. Robust Positive Control Invariant Safe Set:**

A *robust positive control invariant* (RPCI) safe set  $\Omega_{\text{safe}}$  is a set contained in  $\mathbb{X}_c$  such that, if we start within the safe set, there exists a feedback policy  $\pi(\mathbf{x})$  to keep the system in the safe set despite all possible model errors and process noise, captured by  $\mathbb{D}$ .

---

---

**M24. Control Lyapunov Functions:** *Control Lyapunov Functions (CLFs)* extend the notion of Lyapunov stability to continuous-time, control-affine systems (Equation 17). A positive definite function  $L_c$  is a CLF, if there exists a state-feedback control input  $\mathbf{u}$ , such that the time derivative  $\dot{L}_c$  satisfies:  $\dot{L}_c(\mathbf{x}) \leq -L_c(\mathbf{x})$ . The existence of a CLF guarantees that there also exists a state-feedback controller  $\boldsymbol{\pi}(\mathbf{x})$  that asymptotically stabilizes the system.

---

the set  $\Omega_{\text{safe}}$ :

$$\mathbf{u}_{\text{safe},k} = \arg \min_{\mathbf{u}_k \in \mathbb{U}_c} \|\mathbf{u}_k - \mathbf{u}_{\text{learn},k}\|_2^2 \quad 16a.$$

$$\begin{aligned} \text{s.t. } \quad & \mathbf{x}_{k+1} = \bar{\mathbf{f}}_k(\mathbf{x}_k, \mathbf{u}_k) + \hat{\mathbf{f}}_k(\mathbf{x}_k, \mathbf{u}_k, \mathbf{w}_k) \in \Omega_{\text{safe}}, \quad 16b. \\ & \forall \hat{\mathbf{f}}_k(\mathbf{x}_k, \mathbf{u}_k, \mathbf{w}_k) \in \mathbb{D}(\mathbf{x}_k, \mathbf{u}_k), \end{aligned}$$

where the range of possible disturbances  $\mathbb{D}(\mathbf{x}_k, \mathbf{u}_k)$  is given. Since the safety filter and controller are usually decoupled, suboptimal behavior can emerge as the learning-based controller may try to violate the constraints (65).

**Control Barrier Functions (CBFs).** CBFs are used to define safe sets. More specifically, the safe set  $\Omega_{\text{safe}}$  is defined as the superlevel set of a continuously differentiable CBF  $B_c$ ,  $B_c : \mathbb{R}^{n_x} \rightarrow \mathbb{R}$ , as  $\Omega_{\text{safe}} = \{\mathbf{x} \in \mathbb{R}^{n_x} : B_c(\mathbf{x}) \geq 0\}$ . The function  $B_c$  is generally considered for continuous-time, control-affine systems (see also M4) of the form:

$$\dot{\mathbf{x}} = \mathbf{f}_x(\mathbf{x}) + \mathbf{f}_u(\mathbf{x})\mathbf{u}, \quad 17.$$

where  $\mathbf{f}_x$  and  $\mathbf{f}_u$  are locally Lipschitz and  $\mathbf{x}, \mathbf{u}$  are functions of time (118). In the simplest form, the function  $B_c$  is a CBF if there exists a state-feedback control input  $\mathbf{u}$ , such that the time derivative  $\dot{B}_c(\mathbf{x}) = \frac{\partial B_c}{\partial \mathbf{x}} \dot{\mathbf{x}}$  satisfies (118):

$$\sup_{\mathbf{u} \in \mathbb{U}} \frac{\partial B_c}{\partial \mathbf{x}} (\mathbf{f}_x(\mathbf{x}) + \mathbf{f}_u(\mathbf{x})\mathbf{u}) \geq -B_c(\mathbf{x}). \quad 18.$$

The CBF condition in Equation 18 is a continuous-time version of the robust positive control invariance constraint in Equation 16b. In addition to the RPCI constraint, a similar constraint for asymptotic stability can be added to Equation 16 in the form of a constraint on the time derivative of a control Lyapunov function (CLF, see M24)  $L_c$ . However, uncertain dynamics also yield uncertain time derivatives of  $B_c$  and  $L_c$ .

Learning-based approaches extend CBF and CLF analyses to control-affine systems in Equation 17 with known nominal  $\bar{\mathbf{f}}_x, \bar{\mathbf{f}}_u$  and unknown  $\hat{\mathbf{f}}_x, \hat{\mathbf{f}}_u$ . The time derivative of the unknown dynamics for CLFs and/or CBFs (e.g.,  $\frac{\partial B_c}{\partial \mathbf{x}} (\hat{\mathbf{f}}_x(\mathbf{x}) + \hat{\mathbf{f}}_u(\mathbf{x})\mathbf{u})$ ) can be learned from iterative trials (119, 120) or data collected by an RL agent (121, 122). Given a CBF and/or CLF for the true dynamics  $\mathbf{f}$ , improving the estimate of the CBF's or CLF's time derivative for the unknown dynamics  $\hat{\mathbf{f}}$  using data, collected either offline or online, yields a more precise estimate of the constraint in Equation 18. However, any learning error in the CBF's or CLF's time derivative of the unknown dynamics can still lead to applying falsely certified control inputs. To this end, ideas from robust control can be used to guarantee set invariance (Safety Level III) during the learning process by mapping a bounded uncertainty in the dynamics to a bounded uncertainty in the time derivatives of the CBF or CLF (123, 124), or by accounting for all model errors consistent with the collected data (125). In addition, adaptive control approaches have also been proposed to allow safe adaptation of parametric uncertainties in the time derivatives of the CBF or CLF (126, 127). Probabilistic learning techniques for CBFs and CLFs have been used to achieve set invariance probabilistically (Safety Level II) with varying assumptions on the system dynamics: the function  $\mathbf{f}_u(\mathbf{x})$  is fully known (128, 129), a nominal model is known (130), and no nominal model is available (131).

**Hamilton-Jacobi Reachability Analysis.** Another approach for state constraint set certification of a learning-based controller is via the Hamilton-Jacobi (HJ) reachability analysis. The HJ reachability analysis provides a means to estimate a robust positive control invariant safe set  $\Omega_{\text{safe}}$  under dynamics uncertainties (see M23 and Equation 16). Consider a nonlinear system subject to unknown but bounded disturbances  $\hat{\mathbf{f}}(\mathbf{x}) \in \mathbb{D}(\mathbf{x})$ , where  $\mathbb{D}(\mathbf{x})$  is

assumed to be known but possibly conservative. To compute  $\Omega_{\text{safe}}$ , a two-player, zero-sum differential game is formulated:

$$V(\mathbf{x}) = \max_{\mathbf{u}_{\text{sig}} \in \mathbb{U}_{\text{sig}}} \min_{\hat{\mathbf{f}}_{\text{sig}} \in \mathbb{D}_{\text{sig}}} \left( \inf_{k \geq 0} l_c \left( \phi(\mathbf{x}, k; \mathbf{u}_{\text{sig}}, \hat{\mathbf{f}}_{\text{sig}}) \right) \right), \quad 19.$$

where  $V$  is the value function associated with a point  $\mathbf{x} \in \mathbb{X}$ ,  $l_c : \mathbb{X} \mapsto \mathbb{R}$  is a cost function that is non-negative for  $\mathbf{x} \in \mathbb{X}_c$  and negative otherwise,  $\phi(\mathbf{x}, k; \mathbf{u}_{\text{sig}}, \hat{\mathbf{f}}_{\text{sig}})$  denotes the state at  $k$  along a trajectory initialized at  $\mathbf{x}$  following input signal  $\mathbf{u}_{\text{sig}}$  and disturbance signal  $\hat{\mathbf{f}}_{\text{sig}}$ , and  $\mathbb{U}_{\text{sig}}$  and  $\mathbb{D}_{\text{sig}}$  are collections of input and disturbance signals such that each time instance is in  $\mathbb{U}$  and  $\mathbb{D}$ , respectively. The value function  $V$  can be found as the unique viscosity solution of the Hamilton-Jacobi Isaacs (HJI) variational inequality (132). The safe set is then  $\Omega_{\text{safe}} = \{\mathbf{x} \in \mathbb{X} \mid V(\mathbf{x}) \geq 0\}$ . Based on this formulation, we can also obtain an optimally safe policy  $\pi_{\text{safe}}^*$  that maximally steers the system towards the safe set  $\Omega_{\text{safe}}$  (i.e., in the greatest ascent direction of  $V$ ). The HJ reachability analysis allows us to define a safety filter for learning-based control approaches to guarantee constraint set satisfaction (Safety Level II or III). In particular, given  $\Omega_{\text{safe}}$  and  $\pi_{\text{safe}}^*$ , one can safely learn in the interior of  $\Omega_{\text{safe}}$  and apply the optimally safe policy  $\pi_{\text{safe}}^*$  if the system reaches the boundary of  $\Omega_{\text{safe}}$ . To reduce the conservativeness of the approach, a GP-based learning scheme is proposed in (133) to adapt (and shrink) the unknown dynamics set  $\mathbb{D}(\mathbf{x})$  based on observed data.

The general HJ reachability analysis framework (132) has also been combined with on-line dynamics model learning for a target tracking task (134), with online planning for safe exploration (135), and with temporal difference algorithms for safe RL (136). In (137), HJ reachability analysis and CBFs have been integrated to compute smoother control policies while circumventing the need to hand-design appropriate CBFs. Another recent extension (138) proposed modifications that improve the scalability of the HJ safety analysis approach for higher-dimensional systems and was demonstrated on a ten-dimensional quadrotor tracking problem.

**Predictive Safety Filters.** Predictive safety filters can augment any learning-based controller to enforce state constraints,  $\mathbf{x} \in \mathbb{X}_c$ , and input constraints,  $\mathbf{u} \in \mathbb{U}_c$ . They do this by defining the safe invariant set  $\Omega_{\text{safe}}$  in Equation 16b as the set of states (at the next time step) where a sequence of safe control inputs (e.g., from a backup controller) exists that allows the return to (i) a terminal safe set  $\mathbb{X}_{\text{term}}$  (see M25), or (ii) to previously visited safe states.

**Model Predictive Safety Certification (MPSC)** uses the theory of robust MPC in Sec. 2.2 and learning-based robust MPC in Sec. 3.1.3 to filter the output of any learning-based controller, such as of an RL method, to ensure robust constraint satisfaction. The simplest implementation of MPSC (139) uses tube-based MPC and considers the constraints in Equations 12b–12d but replaces the cost in Equation 12a with the cost in Equation 16a to find the closest input  $\mathbf{u}_k$  to the learned input  $\mathbf{u}_{\text{learn},k}$  at the current time step that guarantees that we will continue to satisfy state and input constraints in the future. The main difference between MPSC and learning-based robust MPC described in Sec. 3.1.3 is that the terminal safe set  $\mathbb{X}_{\text{term}}$  in Equation 12d is not coupled with the selection of the cost function in Equation 12a. Instead, the terminal safe set is conservatively initialized with  $\mathbb{X}_{\text{term}} = \Omega_{\text{tube}}$  and can grow to include state trajectories from previous iterations. This approach has been extended to probabilistic constraints by considering a probabilistic tube  $\Omega_{\text{tube}}$  in (140), and to nonlinear nominal models in (141) (Safety Level II).

**Backup control for safe exploration** ensures hard state constraint satisfaction (Safety Level III) by finding a safe backup controller for any given RL policy  $\pi$  (142). Under the assumptions of a known bound  $\mathbb{D}(\mathbf{x}, \mathbf{u})$  on the dynamics  $\mathbf{f}$  and a distance measure to the state constraints  $\mathbb{X}_c$ , the backup controller is used to obtain a future state in the

---

**M25. Terminal Safe Set:** A *terminal safe set*, denoted by  $\mathbb{X}_{\text{term}}$ , is a subset of the safe state constraint space wherein a known auxiliary controller is guaranteed to preserve the system’s state. Entering the terminal safe set at some fixed horizon  $H$  is enforced as a constraint.

---

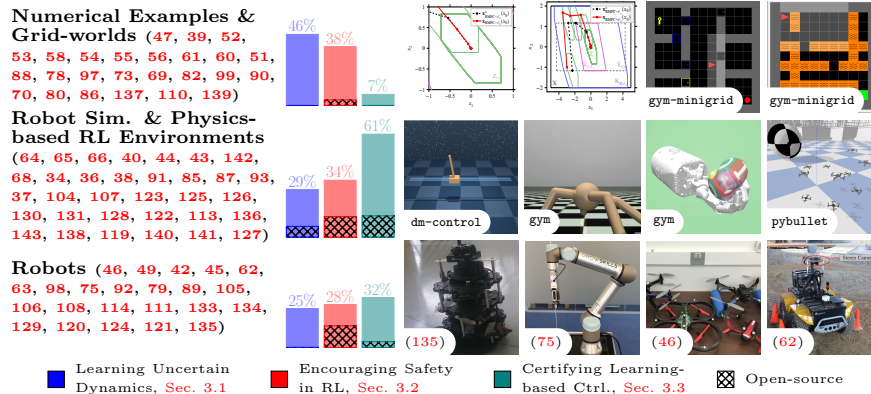


Figure 5: Summary of the environments used for evaluation. With increasing complexity, they can be classified as: abstract numerical examples, robot simulations, and real-world robot experiments. The histograms show the prevalence of each in [Sec. 3.1](#), [Sec. 3.2](#), and [Sec. 3.3](#), as well as the fraction of those whose code is open-sourced.

neighborhood of a previously visited safe state in some prediction horizon. Before a control input  $\mathbf{u}_k$  from  $\pi$  is applied to the system, all possible predicted states  $\mathbf{x}_{k+1}$  must satisfy (i)  $\mathbf{x}_{k+1} \in \mathbb{X}_c$  and (ii) the existence of a safe backup action  $\mathbf{u}_{\text{certified}}(\mathbf{x}_{k+1})$ . Otherwise, the previous backup control input  $\mathbf{u}_{\text{certified}}(\mathbf{x}_k)$  is applied. This procedure guarantees that the system state stays inside a robust positive control invariant set  $\Omega_{\text{safe}} \subseteq \mathbb{X}_c$ .

#### 4. BENCHMARKS

The approaches presented in [Sec. 3](#) have been evaluated in vastly different ways, see [Figure 5](#). The trends we observe are: The works learning uncertain dynamics ([Sec. 3.1](#)) include a preponderance of abstract numerical examples; encouraging safety in RL ([Sec. 3.2](#)) mostly leverages simulations—often based on physics engines ([144](#)) like MuJoCo—but grid worlds are also common; and works certifying learning-based control ([Sec. 3.3](#)) are still mostly simulated but also has the largest fraction of real-world experiments. While numerical examples make it difficult to gauge the practical applicability of a method, we also note that even many RL environments, including physics-based ones, are not representative of existing robotic platforms—for example, they often are deterministic and do not account for variations in the environment.

In an ideal world, all research would be demonstrated in simulations that closely resemble the target system—and brought to (ideally different) real robots, whenever possible. Furthermore, only a minority of the published research open-sourced their software implementations. Even in RL, where sharing code is more common (see, e.g., red bars in [Figure 5](#))—and standard environments and interfaces such as `gym` ([145](#)) have been proposed—the reproducibility of results (that often rely on careful hyper-parameter tuning) remains challenging ([71](#)). With regard to safety, simple RL environments augmented with constraint evaluation ([18](#)) and disturbances ([5](#)) have been proposed, but lack a unified simulation interface for both safe RL and learning-based control approaches—notably, one that also exposes the available *a priori* knowledge of a dynamical system.

We believe that a necessary stepping stone for the advancement of safe learning control is to create physics-based environments that are (i) simple enough to promote adoption, (ii) realistic enough to represent meaningful robotic platforms, (iii) equipped with intuitive interfaces for both control and RL researchers, and (iv) provide useful comparison metrics (e.g., the amount of data required by different approaches).

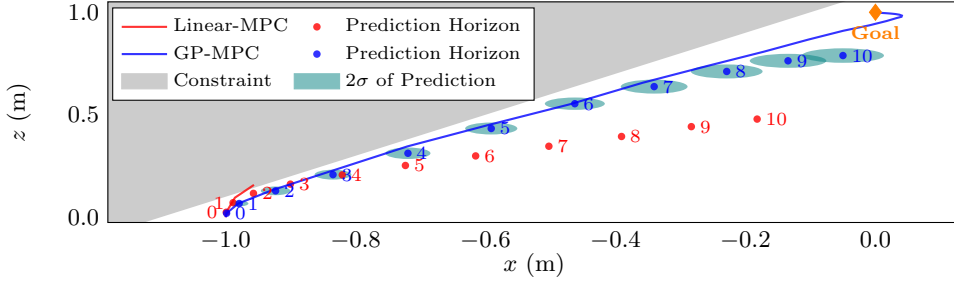


Figure 6: A position comparison of the 2D quadrotor stabilization using linear MPC (red) and GP-MPC (blue), along with the prediction horizons at the second time step, subject to a diagonal state constraint (gray) and input constraints.

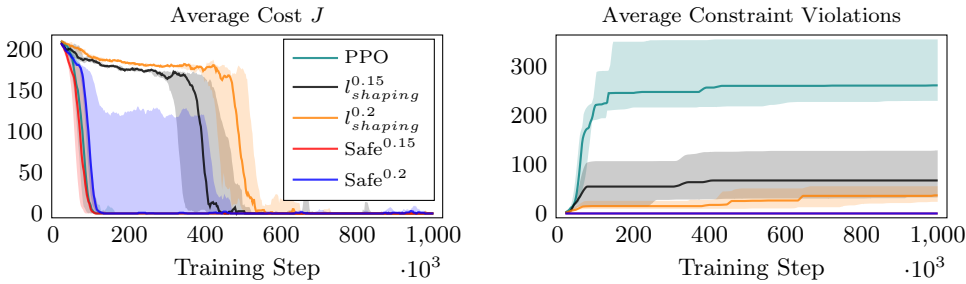


Figure 7: Total cost and constraint violations during learning for PPO, PPO with cost shaping, and PPO with safety layer (safe exploration). Plotted are medians with upper and lower quantiles over 10 seeds. The parameters in superscript represent values of the slack variable, which controls responsiveness to near constraint violations.

**Cart-Pole and Quadrotor Benchmark Environments.** For this reason, we have created an open-source benchmark suite<sup>2</sup> simulating two popular platforms for the evaluation of control and RL research: (i) the cart-pole (64, 104, 136) and (ii) a quadrotor (142, 66, 44, 130, 146). Our simulation environments are based on the open-source Bullet physics engine (147) and we adopt OpenAI’s gym interface (API) (145) for seamless integration with many of the current RL libraries. What sets apart our implementation from previous attempts to create safety-aware RL environments (18, 5) is the extension of the traditional API (145) with features to facilitate: (i) the evaluation of safe learning control approaches (such as the randomization of inertial properties and of initial positions and velocities, random and adversarial disturbances, and the specification and evaluation of state and input constraints) and (ii) the integration with approaches developed by the control theory community that leverage the knowledge of nominal models—notably, we use the symbolic framework CasADi (148) to enrich our gym environments with symbolic *a priori* dynamics, constraints, and cost functions.

**Safe Learning Control Results.** We focus on a constrained stabilization task, in each of our environments (`cartpole` and `quadrotor`). The results we present here are not meant to establish the superiority of one approach over another. Instead, we show how the approaches in (63, 70, 139), taken from each of Sec. 3.1 (learning uncertain dynamics), Sec. 3.2 (encouraging safety in RL), and Sec. 3.3 (certifying learning-based control), can improve

<sup>2</sup>Safe control benchmark suite on GitHub: <https://github.com/utiasDSL/safe-control-gym>

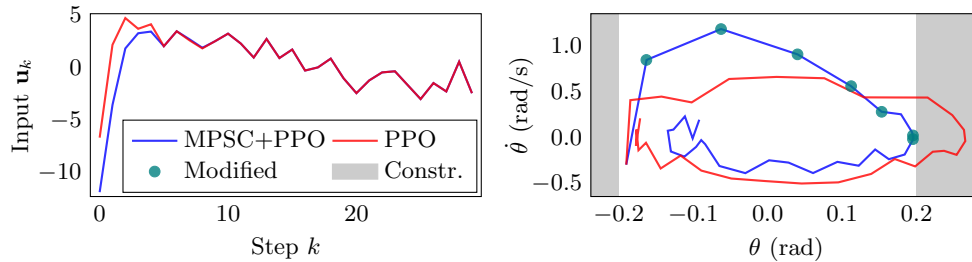


Figure 8: Left: The uncertified PPO input (red) is plotted against the certified MPSC+PPO input (blue). Right: The cart-pole state diagram ( $\theta$  and  $\dot{\theta}$ ) comparing the MPSC+PPO certified trajectory (blue) and the uncertified PPO trajectory (red). The MPSC is most active (green dots) when the system is about to leave the constraint boundary (gray) or the set of states from which the MPSC can correct the system.

control performance while pursuing constraint satisfaction. By doing so, we demonstrate that our benchmark can be easily integrated with learnable policies developed by either the control or RL research community. This also allows us to fairly compare the data hungeriness of the different safe learning control approaches.

A learning-based robust MPC with a GP estimate of  $\hat{f}$  (GP-MPC), representative of [Sec. 3.1](#) (learning uncertain dynamics), was implemented to stabilize a quadrotor subject to a state and input constraints, following the approach in [\(63\)](#). For this experiment, a linearization about hover with a mass and moment of inertia 150% of the true values was used as the prior model. Eight-hundred randomly selected state-action pairs (or 80 seconds), sampled from a couple of minutes of training data, were used for hyper-parameter optimization, which was performed offline. [Figure 6](#) compares the performance of a linear MPC, using the inaccurate prior model, with the GP-MPC approach. We see that the Linear MPC, using the heavier prior model, predicts the trajectory of the quadrotor will be relatively shallow when maximum thrust is applied. This results in the quadrotor quickly violating the position constraint. In contrast, the GP-MPC is able to account for the inaccurate prior model and respects the constraint boundary by a margin proportional to the 95% confidence interval on its dynamics prediction, stabilizing the quadrotor.

We also applied the Safe Exploration approach [\(70\)](#) to the popular deep RL algorithm PPO [\(149\)](#), as a representative of [Sec. 3.2](#) (encouraging safety in RL), and tested it on the cart-pole stabilization task with constraints on the cart position. Notably, the task terminates upon any constraint violation. We compared it against two baselines, standard PPO and PPO with naive cost shaping (adding a penalty when close to constraint violation). Each RL approach used over 9 hours of simulation time collecting data for training. [Figure 7](#) shows that the two constraint-aware approaches (cost shaping and safe exploration [\(70\)](#)) achieve the same performance after learning, but with substantially lower constraint violations than the standard PPO. The safe exploration approach outperforms cost shaping in terms of constraint satisfaction, while not compromising convergence speed towards reaching the optimal cost. Safe Exploration, however, requires careful parameter tuning, especially on the slack variable that determines its responsiveness to near constraint violation.

Finally, a Model Predictive Safety Certification (MPSC) algorithm, based on [\(139\)](#), was chosen as a representative of [Sec. 3.3](#) (certifying learning-based control). This particular formulation uses an MPC framework to modify an unsafe learning controller’s actions. Here, a sub-optimal PPO controller provides the uncertified inputs trying to stabilize the cart-pole. The advantages of using MPSC are highlighted in [Figure 8](#). In [Figure 8](#) (right), the inputs are modified by the MPSC early in the stabilization to keep the cart-pole within the



constraint boundaries. **Figure 8** (left) shows that, without the MPSC, PPO violates the constraints, but with MPSC, it manages to stay within the boundaries. The plot also shows that MPSC is most active when the system is close to the constraint boundaries (the green dots show when MPSC modified the learning controller’s input). This provides a proof of concept of how safety filters can be combined with RL control for safer performance.

In our future work, we intend to use our `cartpole` and `quadrotor` benchmark environments to evaluate the robustness, performance, safety and data efficiency of different control and learning approaches.

## 5. DISCUSSION AND PERSPECTIVES FOR FUTURE DIRECTIONS

The problem of safe learning control is emerging as a crucial topic for next-generation robotics. In this review, we summarized approaches from the control and the machine learning communities that allow data to be safely used to improve the closed-loop performance of robot control systems. We show that machine learning techniques, particularly RL, can help generalization towards larger classes of systems (i.e., fewer prior model assumptions), while control theory provides the insights and frameworks necessary to provide constraint satisfaction guarantees and closed-loop stability guarantees during the learning. Despite the many advances to date, there remain many opportunities for future research.

### OPEN CHALLENGES

1. **Capturing a Broader Class of Systems.** Work to date has focused on nonlinear systems in the form of **Equation 1**. While they can model many robotic platforms, robots can also exhibit hybrid dynamics (e.g., legged robots or other contact dynamics with the environment (150)), time-varying dynamics (e.g., operation in changing environments (151, 152)), time delays (e.g., in actuation, sensing, or observing the reward (6)), partial differential dependencies in the dynamics (e.g., in continuum robotics (153)). Expanding safe learning control approaches to these scenarios is essential for their broader applicability in robotics (see **Figure 4**).
2. **Accounting for Imperfect State Measurements.** The majority of safe learning control approaches assume direct access to (possibly noisy) state measurements and neglect the problem of state estimation. In practice, obtaining accurate state information is challenging due to sensors that do not provide state measurements directly (e.g., images as measurements), inaccurate process and observation models used for state estimation, and/or improper state feature representations. One open challenge is to account for state estimation errors and learned process and observation models in safe learning control (143). Expanding existing approaches to work with (possibly high-dimensional) sensor data is essential for a broad applicability of these methods in robotics.
3. **Considering Scalability, and Sampling and Computational Efficiency.** Many of the approaches presented here have only been demonstrated on small toy problems and applying them to high-dimensional robotics problems is not trivial. Moreover, in practice, we often face issues such as data sparsity, distribution shifts, and the optimality-complexity trade-off for real-time implementations. Efficient robot learning relies on multiple factors including control architecture design (154), systematic training data collection (155), and appropriate function class selection (156). While current approaches focus on providing theoretical safety guarantees, formal analysis of sampling complexity and computational complexity is indispensable to facilitate the implementation of safe learning control algorithms in real-world robot applications.

4. **Verifying System and Modeling Assumptions.** The safety guarantees provided often rely on a set of assumptions (e.g., Lipschitz continuous *true* dynamics with a known Lipschitz constant or bounded disturbance sets). It is difficult to verify these assumptions prior to a robot’s operation. To facilitate algorithm implementation, we also see other approximations being made (e.g., linearization, data assumed to be independent and identically distributed (i.i.d.) Gaussian samples). Systematic approaches to verify or quantify the impact of the assumptions and the approximations with minimal (online) data are crucial to allow the safe learning approaches to be applied in real-world applications. This can also include investigations into the interpretability of trained models, especially black-box models such as DNNs, for safe closed-loop operation (6).

This list of challenges is by no means complete. Other important and open questions pertinent to safe learning control are:

- What are the appropriate **benchmarks, evaluation metrics, and practices** to provide practical insights and perform fair comparisons of algorithms that rely on different assumptions?
- What should be the **role of simulation** in the offline design and evaluation phases? Can simulation be used to find safe hyperparameters?
- How do we define safety in **human-robot interaction**?
- How can data be safely used in **multi-agent learning** settings?

Efforts that combine control theory and machine learning for safe learning control have been shown to result in improved control performance and system safety. Recent successes and growing interest should motivate the further development of a systematic body of theory, advanced methodologies, and computational methods for safe learning in robotics, bringing the large potential of learning into safety-critical control and robotics applications. The open challenges and questions are intended to push us further, towards a future of real-world safe autonomy where robots reliably and safely function in complex environments.

## ACKNOWLEDGMENTS

The authors would like to acknowledge the early contributions to this work by Karime Pereira and Sepehr Samavi, the invaluable suggestions and feedback by Hallie Siegel, as well as the support from the Natural Sciences and Engineering Research Council of Canada (NSERC), the Canada Research Chairs Program, and the CIFAR AI Chair.

## LITERATURE CITED

1. Burnett K, Qian J, Du X, Liu L, Yoon DJ, et al. 2021. Zeus: A system description of the two-time winner of the collegiate SAE autodrive competition. *Journal of Field Robotics* 38:139–166
2. Boutilier JJ, Brooks SC, Janmohamed A, Byers A, Buick JE, et al. 2017. Optimizing a drone network to deliver automated external defibrillators. *Circulation* 135:2454–2465
3. Dong K, Pereira K, Shkurti F, Schoellig AP. 2020. Catch the ball: Accurate high-speed motions for mobile manipulators via inverse dynamics learning. [arXiv:2003.07489](#) [cs.RO]
4. García J, Fern, o Fernández. 2015. A comprehensive survey on safe reinforcement learning. *Journal of Machine Learning Research* 16:1437–1480
5. Dulac-Arnold G, Levine N, Mankowitz DJ, Li J, Paduraru C, et al. 2021. An empirical investigation of the challenges of real-world reinforcement learning. [arXiv:2003.11881](#) [cs.LG]
6. Dulac-Arnold G, Mankowitz D, Hester T. 2019. Challenges of real-world reinforcement learning. [arXiv:1904.12901](#) [cs.LG]

7. Hewing L, Wabersich KP, Menner M, Zeilinger MN. 2020. Learning-based model predictive control: Toward safe learning in control. *Annual Review of Control, Robotics, and Autonomous Systems* 3:269–296
8. Bristow D, Tharayil M, Alleyne A. 2006. A survey of iterative learning control. *IEEE Control Systems Magazine* 26:96–114
9. Ahn HS, Chen Y, Moore KL. 2007. Iterative learning control: Brief survey and categorization. *IEEE Transactions on Systems, Man, and Cybernetics, Part C (Applications and Reviews)* 37:1099–1121
10. Polydoros AS, Nalpantidis L. 2017. Survey of model-based reinforcement learning: Applications on robotics. *Journal of Intelligent & Robotic Systems* 86:153–173
11. Chazilygeroudis K, Vassiliades V, Stulp F, Calinon S, Mouret JB. 2020. A survey on policy search algorithms for learning robot controllers in a handful of trials. *IEEE Transactions on Robotics* 36:328–347
12. Ravichandar H, Polydoros AS, Chernova S, Billard A. 2020. Recent advances in robot learning from demonstration. *Annual Review of Control, Robotics, and Autonomous Systems* 3:297–330
13. Kober J, Bagnell JA, Peters J. 2013. Reinforcement learning in robotics: A survey. *The International Journal of Robotics Research* 32:1238–1274
14. Recht B. 2019. A tour of reinforcement learning: The view from continuous control. *Annual Review of Control, Robotics, and Autonomous Systems* 2:253–279
15. Kiumarsi B, Vamvoudakis KG, Modares H, Lewis FL. 2018. Optimal and autonomous control using reinforcement learning: A survey. *IEEE Transactions on Neural Networks and Learning Systems* 29:2042–2062
16. Osborne M, Shin HS, Tsourdos A. 2021. *A Review of Safe Online Learning for Nonlinear Control Systems*. In *2021 International Conference on Unmanned Aircraft Systems (ICUAS)*, pp. 794–803. Piscataway, NJ: IEEE
17. Tambon F, Laberge G, An L, Nikanjam A, Mindom PSN, et al. 2021. How to certify machine learning based safety-critical systems? a systematic literature review. [arXiv:2107.12045](https://arxiv.org/abs/2107.12045) [cs.LG]
18. Ray A, Achiam J, Amodei D. 2019. Benchmarking Safe Exploration in Deep Reinforcement Learning. <https://cdn.openai.com/safexp-short.pdf>
19. Leike J, Martic M, Krakovna V, Ortega PA, Everitt T, et al. 2017. AI safety gridworlds. [arXiv:1711.09883](https://arxiv.org/abs/1711.09883) [cs.LG]
20. Sutton RS, Barto AG. 2018. *Reinforcement Learning: An Introduction*. Cambridge, MA: MIT Press, 2nd ed.
21. Khalil H. 2002. *Nonlinear Systems*. Pearson Education. Prentice Hall
22. Åström K, Wittenmark B. 2011. *Computer-Controlled Systems: Theory and Design, Third Edition*. Dover Books on Electrical Engineering. Dover Publications
23. Sastry S, Bodson M. 2011. *Adaptive Control: Stability, Convergence and Robustness*. Dover Books on Electrical Engineering Series. Dover Publications
24. Nguyen-Tuong D, Peters J. 2011. Model learning for robot control: a survey. *Cognitive processing* 12:319–340
25. Zhou K, Doyle J, Glover K. 1996. *Robust and Optimal Control*. Prentice Hall
26. Dullerud G, Paganini F. 2005. *A Course in Robust Control Theory: A Convex Approach*. Texts in Applied Mathematics. Springer New York
27. Rawlings J, Mayne D, Diehl M. 2017. *Model Predictive Control: Theory, Computation, and Design*. Nob Hill Publishing
28. Mayne D, Seron M, Raković S. 2005. Robust model predictive control of constrained linear systems with bounded disturbances. *Automatica* 41:219–224
29. Arulkumaran K, Deisenroth MP, Brundage M, Bharath AA. 2017. Deep reinforcement learning: A brief survey. *IEEE Signal Processing Magazine* 34:26–38
30. Dai B, Shaw A, Li L, Xiao L, He N, et al. 2018. *SBEED: Convergent Reinforcement Learning with Nonlinear Function Approximation*. In *Proceedings of the 35th International Conference on Machine Learning*, ed. J Dy, A Krause, pp. 1125–1134, vol. 80 of *Proceedings of Machine Learning Research*, pp. 1125–1134. N.p.: PMLR
31. Cheng R, Verma A, Orosz G, Chaudhuri S, Yue Y, Burdick J. 2019. Control regularization for reduced variance reinforcement learning. In *Proceedings of the 36th International Conference on Machine Learning*, ed. K Chaudhuri, R Salakhutdinov, pp. 1141–1150, vol. 97 of

- Proceedings of Machine Learning Research*. N.p.: PMLR
32. Ghavamzadeh M, Mannor S, Pineau J, Tamar A. 2015. Bayesian reinforcement learning: A survey. *Foundations and Trends® in Machine Learning* 8:359–483
  33. Altman E. 1999. *Constrained Markov decision processes*, vol. 7. Boca Raton, FL: Chapman & Hall/CRC Press
  34. Achiam J, Held D, Tamar A, Abbeel P. 2017. Constrained policy optimization. In *Proceedings of the 34th International Conference on Machine Learning*, ed. D Precup, YW Teh, pp. 22–31, vol. 70 of *Proceedings of Machine Learning Research*. N.p.: PMLR
  35. Nilim A, El Ghaoui L. 2005. Robust control of Markov decision processes with uncertain transition matrices. *Operations Research* 53:780–798
  36. Pinto L, Davidson J, Sukthankar R, Gupta A. 2017. Robust adversarial reinforcement learning. In *Proceedings of the 34th International Conference on Machine Learning*, ed. D Precup, YW Teh, pp. 2817–2826, vol. 70 of *Proceedings of Machine Learning Research*. N.p.: PMLR
  37. Pan X, Seita D, Gao Y, Canny J. 2019. *Risk Averse Robust Adversarial Reinforcement Learning*. In *2019 International Conference on Robotics and Automation (ICRA)*, pp. 8522–8528. Piscataway, NJ: IEEE
  38. Vinitzky E, Du Y, Parvate K, Jang K, Abbeel P, Bayen A. 2020. Robust reinforcement learning using adversarial populations. [arXiv:2008.01825](https://arxiv.org/abs/2008.01825) [cs.LG]
  39. Cooper J, Che J, Cao C. 2014. The use of learning in fast adaptation algorithms. *International Journal of Adaptive Control and Signal Processing* 28:325–340
  40. Gahlawat A, Zhao P, Patterson A, Hovakimyan N, Theodorou E. 2020. L1-GP: L1 adaptive control with Bayesian learning. In *Proceedings of the 2nd Conference on Learning for Dynamics and Control*, ed. AM Bayen, A Jadbabaie, G Pappas, PA Parrilo, B Recht, C Tomlin, M Zeilinger, pp. 826–837, vol. 120 of *Proceedings of Machine Learning Research*. N.p.: PMLR
  41. Hovakimyan N, Cao C. 2010. *L1 Adaptive Control Theory: Guaranteed Robustness with Fast Adaptation*. USA: Society for Industrial and Applied Mathematics
  42. Grande RC, Chowdhary G, How JP. 2014. Experimental validation of Bayesian nonparametric adaptive control using Gaussian processes. *Journal of Aerospace Information Systems* 11:565–578
  43. Chowdhary G, Kingravi HA, How JP, Vela PA. 2015. Bayesian nonparametric adaptive control using Gaussian processes. *IEEE Transactions on Neural Networks and Learning Systems* 26:537–550
  44. Joshi G, Chowdhary G. 2019. *Deep model reference adaptive control*. In *2019 IEEE 58th Conference on Decision and Control (CDC)*, pp. 4601–4608. Piscataway, NJ: IEEE
  45. Joshi G, Virdi J, Chowdhary G. 2020. Asynchronous deep model reference adaptive control. [arXiv:2011.02920](https://arxiv.org/abs/2011.02920) [cs.RO]
  46. Berkenkamp F, Schoellig AP. 2015. *Safe and robust learning control with Gaussian processes*. In *2015 European Control Conference (ECC)*, pp. 2496–2501. Piscataway, NJ: IEEE
  47. Holicki T, Scherer CW, Trimpe S. 2021. Controller design via experimental exploration with robustness guarantees. *IEEE Control Systems Letters* 5:641–646
  48. von Rohr A, Neumann-Brosig M, Trimpe S. 2021. Probabilistic robust linear quadratic regulators with Gaussian processes. In *Proceedings of the 3rd Conference on Learning for Dynamics and Control*, ed. A Jadbabaie, J Lygeros, GJ Pappas, PA Parrilo, B Recht, CJ Tomlin, MN Zeilinger, pp. 324–335, vol. 144 of *Proceedings of Machine Learning Research*. N.p.: PMLR
  49. Helwa MK, Heins A, Schoellig AP. 2019. Provably robust learning-based approach for high-accuracy tracking control of lagrangian systems. *IEEE Robotics and Automation Letters* 4:1587–1594
  50. Greeff M, Schoellig AP. 2021. Exploiting differential flatness for robust learning-based tracking control using Gaussian processes. *IEEE Control Systems Letters* 5:1121–1126
  51. Tanaskovic M, Fagiano L, Smith R, Morari M. 2014. Adaptive receding horizon control for constrained MIMO systems. *Automatica* 50:3019–3029
  52. Lorenzen M, Cannon M, Allgöwer F. 2019. Robust MPC with recursive model update. *Automatica* 103:461–471
  53. Bujarbaruah M, Zhang X, Borrelli F. 2018. *Adaptive MPC with Chance Constraints for FIR Systems*. In *2018 Annual American Control Conference (ACC)*, pp. 2312–2317. Piscataway, NJ: IEEE
  54. Bujarbaruah M, Zhang X, Tanaskovic M, Borrelli F. 2019. Adaptive MPC under Time Varying

- Uncertainty: Robust and Stochastic. [arXiv:1909.13473](https://arxiv.org/abs/1909.13473) [eess.SY]
55. Gonçalves GA, Guay M. 2016. Robust discrete-time set-based adaptive predictive control for nonlinear systems. *Journal of Process Control* 39:111–122
  56. Köhler J, Kötting P, Soloperto R, Allgöwer F, Müller MA. 2020. A robust adaptive model predictive control framework for nonlinear uncertain systems. *International Journal of Robust and Nonlinear Control* <https://doi.org/10.1002/rnc.5147>
  57. Rosolia U, Borrelli F. 2018. Learning model predictive control for iterative tasks. a data-driven control framework. *IEEE Transactions on Automatic Control* 63:1883–1896
  58. Bujarbaruah M, Zhang X, Rosolia U, Borrelli F. 2018. *Adaptive MPC for Iterative Tasks*. In *2018 IEEE Conference on Decision and Control (CDC)*, pp. 6322–6327. Piscataway, NJ: IEEE
  59. Pereida K, Brunke L, Schoellig AP. 2021. Robust adaptive model predictive control for guaranteed fast and accurate stabilization in the presence of model errors. *International Journal of Robust and Nonlinear Control* in-press
  60. Aswani A, Gonzalez H, Sastry SS, Tomlin C. 2013. Provably safe and robust learning-based model predictive control. *Automatica* 49:1216–1226
  61. Soloperto R, Müller MA, Trimpe S, Allgöwer F. 2018. Learning-based robust model predictive control with state-dependent uncertainty. *IFAC-PapersOnLine* 51(20):442–447
  62. Ostafew CJ, Schoellig AP, Barfoot TD. 2016. Robust constrained learning-based NMPC enabling reliable mobile robot path tracking. *The International Journal of Robotics Research* 35:1547–1563
  63. Hewing L, Kabzan J, Zeilinger MN. 2020. Cautious model predictive control using Gaussian process regression. *IEEE Transactions on Control Systems Technology* 28:2736–2743
  64. Kamthe S, Deisenroth M. 2018. Data-efficient reinforcement learning with probabilistic model predictive control. In *Proceedings of the Twenty-First International Conference on Artificial Intelligence and Statistics*, ed. A Storkey, F Perez-Cruz, pp. 1701–1710, vol. 84 of *Proceedings of Machine Learning Research*. N.p.: PMLR
  65. Koller T, Berkenkamp F, Turchetta M, Boedecker J, Krause A. 2019. Learning-based model predictive control for safe exploration and reinforcement learning. [arXiv:1906.12189](https://arxiv.org/abs/1906.12189) [eess.SY]
  66. Fan D, Agha A, Theodorou E. 2020. *Deep Learning Tubes for Tube MPC*. In *Proceedings of Robotics: Science and Systems XVI*, pp. pap. 87. N.p.: Robot. Sci. Syst. Found.
  67. McKinnon CD, Schoellig AP. 2020. *Context-aware Cost Shaping to Reduce the Impact of Model Error in Receding Horizon Control*. In *2020 IEEE International Conference on Robotics and Automation (ICRA)*, pp. 2386–2392. Piscataway, NJ: IEEE
  68. Berkenkamp F, Turchetta M, Schoellig A, Krause A. 2017. Safe model-based reinforcement learning with stability guarantees. In *Advances in Neural Information Processing Systems*, ed. I Guyon, UV Luxburg, S Bengio, H Wallach, R Fergus, S Vishwanathan, R Garnett, pp. 908–919, vol. 30. Red Hook, NY, USA: Curran Associates, Inc.
  69. Turchetta M, Berkenkamp F, Krause A. 2016. Safe exploration in finite Markov decision processes with Gaussian processes. In *Advances in Neural Information Processing Systems*, ed. DD Lee, M Sugiyama, UV Luxburg, I Guyon, R Garnett, pp. 4312–4320, vol. 29. Red Hook, NY, USA: Curran Associates, Inc.
  70. Dalal G, Dvijotham K, Vecerik M, Hester T, Paduraru C, Tassa Y. 2018. Safe exploration in continuous action spaces. [arXiv:1801.08757](https://arxiv.org/abs/1801.08757) [cs.AI]
  71. Henderson P, Islam R, Bachman P, Pineau J, Precup D, Meger D. 2018. *Deep Reinforcement Learning That Matters*. In *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 32(1). Palo Alto, CA: AAAI Press. <https://ojs.aaai.org/index.php/AAAI/article/view/11694>
  72. Mezić I. 2003. Controllability, integrability and ergodicity. In *Multidisciplinary Research in Control: The Mohammed Dahleh Symposium 2002*, ed. L Giarré, B Bamieh, pp. 213–229. Berlin, Heidelberg: Springer Berlin Heidelberg
  73. Moldovan TM, Abbeel P. 2012. *Safe Exploration in Markov Decision Processes*. In *Proceedings of the 29th International Conference on Machine Learning (ICML)*, pp. 1451–1458. Madison, WI: Omnipress
  74. Brafman RI, Tennenholtz M. 2002. R-max—a general polynomial time algorithm for near-optimal reinforcement learning. *Journal of Machine Learning Research* 3:213–231
  75. Pham TH, De Magistris G, Tachibana R. 2018. *OptLayer - Practical Constrained Optimization for Deep Reinforcement Learning in the Real World*. In *2018 IEEE International Conference*

- on *Robotics and Automation (ICRA)*, pp. 6236–6243. Piscataway, NJ: IEEE
76. Kim Y, Allmendinger R, López-Ibáñez M. 2021. Safe learning and optimization techniques: Towards a survey of the state of the art. [arXiv:2101.09505](#) [cs.LG]
  77. Duivenvoorden RR, Berkenkamp F, Carion N, Krause A, Schoellig AP. 2017. Constrained Bayesian optimization with particle swarms for safe adaptive controller tuning. *IFAC-PapersOnLine* 50(1):11800–11807
  78. Sui Y, Gotovos A, Burdick J, Krause A. 2015. Safe exploration for optimization with Gaussian processes. In *Proceedings of the 32nd International Conference on Machine Learning*, ed. F Bach, D Blei, pp. 997–1005, vol. 37 of *Proceedings of Machine Learning Research*. N.p.: PMLR
  79. Berkenkamp F, Krause A, Schoellig AP. 2020. Bayesian optimization with safety constraints: Safe and automatic parameter tuning in robotics. [arXiv:1602.04450](#) [cs.RO]
  80. Sui Y, Vincent Zhuang, Burdick J, Yue Y. 2018. Stagewise safe Bayesian optimization with Gaussian processes. In *Proceedings of the 35th International Conference on Machine Learning*, ed. J Dy, A Krause, pp. 4781–4789, vol. 80 of *Proceedings of Machine Learning Research*. N.p.: PMLR
  81. Baumann D, Marco A, Turchetta M, Trimpe S. 2021. GoSafe: Globally optimal safe robot learning. [arXiv:2105.13281](#) [cs.RO]
  82. Wachi A, Sui Y, Yue Y, Ono M. 2018. *Safe Exploration and Optimization of Constrained MDPs Using Gaussian Processes*. In *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 32(1). Palo Alto, CA: AAAI Press. <https://ojs.aaai.org/index.php/AAAI/article/view/12103>
  83. Kumar A, Zhou A, Tucker G, Levine S. 2020. Conservative Q-learning for offline reinforcement learning. [arXiv:2006.04779](#) [cs.LG]
  84. Chua K, Calandra R, McAllister R, Levine S. 2018. Deep reinforcement learning in a handful of trials using probabilistic dynamics models. In *Advances in Neural Information Processing Systems*, ed. S Bengio, H Wallach, H Larochelle, K Grauman, N Cesa-Bianchi, R Garnett, pp. 4759–4770, vol. 31. Red Hook, NY, USA: Curran Associates, Inc.
  85. Srinivasan K, Eysenbach B, Ha S, Tan J, Finn C. 2020. Learning to be safe: Deep RL with a safety critic. [arXiv:2010.14603](#) [cs.LG]
  86. Thananjeyan B, Balakrishna A, Nair S, Luo M, Srinivasan K, et al. 2021. Recovery RL: Safe reinforcement learning with learned recovery zones. *IEEE Robotics and Automation Letters* 6:4915–4922
  87. Bharadhwaj H, Kumar A, Rhinehart N, Levine S, Shkurti F, Garg A. 2021. Conservative safety critics for exploration. [arXiv:2010.14497](#) [cs.LG]
  88. Chow Y, Ghavamzadeh M, Janson L, Pavone M. 2017. Risk-constrained reinforcement learning with percentile risk criteria. *Journal of Machine Learning Research* 18:1–51
  89. Kahn G, Villaflor A, Pong V, Abbeel P, Levine S. 2017. Uncertainty-aware reinforcement learning for collision avoidance. [arXiv:1702.01182](#) [cs.LG]
  90. Lütjens B, Everett M, How JP. 2019. *Safe Reinforcement Learning With Model Uncertainty Estimates*. In *2019 International Conference on Robotics and Automation (ICRA)*, pp. 8662–8668. Piscataway, NJ: IEEE
  91. Zhang J, Cheung B, Finn C, Levine S, Jayaraman D. 2020. Cautious adaptation for reinforcement learning in safety-critical settings. In *Proceedings of the 37th International Conference on Machine Learning*, ed. HD III, A Singh, pp. 11055–11065, vol. 119 of *Proceedings of Machine Learning Research*. N.p.: PMLR
  92. Thananjeyan B, Balakrishna A, Rosolia U, Li F, McAllister R, et al. 2020. Safety augmented value estimation from demonstrations (SAVED): Safe deep model-based rl for sparse cost robotic tasks. *IEEE Robotics and Automation Letters* 5:3612–3619
  93. Urpí NA, Curi S, Krause A. 2021. Risk-averse offline reinforcement learning. [arXiv:2102.05371](#) [cs.LG]
  94. Bellemare MG, Dabney W, Munos R. 2017. *A Distributional Perspective on Reinforcement Learning*. In *Proceedings of the 34th International Conference on Machine Learning*, ed. D Precup, YW Teh, pp. 449–458, vol. 70 of *Proceedings of Machine Learning Research*, pp. 449–458. N.p.: PMLR
  95. Liang Q, Que F, Modiano E. 2018. Accelerated primal-dual policy optimization for safe reinforcement learning. [arXiv:1802.06480](#) [cs.AI]

96. Schulman J, Levine S, Abbeel P, Jordan M, Moritz P. 2015. Trust region policy optimization. In *Proceedings of the 32nd International Conference on Machine Learning*, ed. F Bach, D Blei, pp. 1889–1897, vol. 37 of *Proceedings of Machine Learning Research*. N.p.: PMLR
97. Chow Y, Nachum O, Duenez-Guzman E, Ghavamzadeh M. 2018. A Lyapunov-based approach to safe reinforcement learning. In *Advances in Neural Information Processing Systems*, ed. S Bengio, H Wallach, H Larochelle, K Grauman, N Cesa-Bianchi, R Garnett, pp. 8103–8112, vol. 31. Red Hook, NY, USA: Curran Associates, Inc.
98. Chow Y, Nachum O, Faust A, Duenez-Guzman E, Ghavamzadeh M. 2019. Lyapunov-based safe policy optimization for continuous control. [arXiv:1901.10031](https://arxiv.org/abs/1901.10031) [cs.LG]
99. Satija H, Amortila P, Pineau J. 2020. Constrained Markov decision processes via backward value functions. In *Proceedings of the 37th International Conference on Machine Learning*, ed. HD III, A Singh, pp. 8502–8511, vol. 119 of *Proceedings of Machine Learning Research*. N.p.: PMLR
100. Morimoto J, Doya K. 2005. Robust reinforcement learning. *Neural Computation* 17:335–359
101. Turchetta M, Krause A, Trimpe S. 2020. *Robust Model-free Reinforcement Learning with Multi-objective Bayesian Optimization*. In *2020 IEEE International Conference on Robotics and Automation (ICRA)*, pp. 10702–10708. Piscataway, NJ: IEEE
102. Goodfellow IJ, Pouget-Abadie J, Mirza M, Xu B, Warde-Farley D, et al. 2014. Generative adversarial networks. [arXiv:1406.2661](https://arxiv.org/abs/1406.2661) [stat.ML]
103. Mnih V, Kavukcuoglu K, Silver D, Rusu AA, Veness J, et al. 2015. Human-level control through deep reinforcement learning. *Nature* 518:529–533
104. Lütjens B, Everett M, How JP. 2020. Certified adversarial robustness for deep reinforcement learning. In *Proceedings of the Conference on Robot Learning*, ed. LP Kaelbling, D Kragic, K Sugiura, pp. 1328–1337, vol. 100 of *Proceedings of Machine Learning Research*. N.p.: PMLR
105. Sadeghi F, Levine S. 2017. *CAD2RL: Real Single-Image Flight Without a Single Real Image*. In *Proceedings of Robotics: Science and Systems XIII*, pp. pap. 34. N.p.: Robot. Sci. Syst. Found.
106. Loquercio A, Kaufmann E, Ranftl R, Dosovitskiy A, Koltun V, Scaramuzza D. 2020. Deep drone racing: From simulation to reality with domain randomization. *IEEE Transactions on Robotics* 36:1–14
107. Rajeswaran A, Ghotra S, Ravindran B, Levine S. 2017. EPOpt: Learning robust neural network policies using model ensembles. [arXiv:1610.01283](https://arxiv.org/abs/1610.01283) [cs.LG]
108. Mehta B, Diaz M, Golemo F, Pal CJ, Paull L. 2020. Active domain randomization. In *Proceedings of the Conference on Robot Learning*, ed. LP Kaelbling, D Kragic, K Sugiura, pp. 1162–1176, vol. 100 of *Proceedings of Machine Learning Research*. N.p.: PMLR
109. Zhou S, Helwa MK, Schoellig AP. 2020. Deep neural networks as add-on modules for enhancing robot performance in impromptu trajectory tracking. *The International Journal of Robotics Research* 39:1397–1418
110. Jin M, Lavaei J. 2020. Stability-certified reinforcement learning: A control-theoretic perspective. *IEEE Access* 8:229086–229100
111. Shi G, Shi X, O’Connell M, Yu R, Azizzadenesheli K, et al. 2019. *Neural Lander: Stable Drone Landing Control Using Learned Dynamics*. In *2019 International Conference on Robotics and Automation (ICRA)*, pp. 9784–9790. Piscataway, NJ: IEEE
112. Fazlyab M, Robey A, Hassani H, Morari M, Pappas GJ. 2019. Efficient and accurate estimation of Lipschitz constants for deep neural networks. [arXiv:1906.04893](https://arxiv.org/abs/1906.04893) [cs.LG]
113. Richards SM, Berkenkamp F, Krause A. 2018. The Lyapunov neural network: Adaptive stability certification for safe learning of dynamical systems. In *Proceedings of The 2nd Conference on Robot Learning*, ed. A Billard, A Dragan, J Peters, J Morimoto, pp. 466–476, vol. 87 of *Proceedings of Machine Learning Research*. N.p.: PMLR
114. Zhou Z, Oguz OS, Leibold M, Buss M. 2020. A general framework to increase safety of learning algorithms for dynamical systems based on region of attraction estimation. *IEEE Transactions on Robotics* 36:1472–1490
115. Jarvis-Wloszek Z, Feeley R, Tan W, Sun K, Packard A. 2003. *Some controls applications of sum of squares programming*. In *42nd IEEE International Conference on Decision and Control (CDC)*, vol. 5, pp. 4676–4681 Vol.5. Piscataway, NJ: IEEE
116. Schilders WH, Van der Vorst HA, Rommes J. 2008. *Model order reduction: theory, research aspects and applications*, vol. 13. Springer-Verlag Berlin Heidelberg

117. Alshiekh M, Bloem R, Udiger Ehlers R, Könighofer B, Niekum S, Topcu U. 2018. *Safe Reinforcement Learning via Shielding*. In *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 32(1). Palo Alto, CA: AAAI Press. <https://ojs.aaai.org/index.php/AAAI/article/view/11797>
118. Ames AD, Coogan S, Egerstedt M, Notomista G, Sreenath K, Tabuada P. 2019. *Control Barrier Functions: Theory and Applications*. In *2019 18th European Control Conference (ECC)*, pp. 3420–3431. Piscataway, NJ: IEEE
119. Taylor AJ, Dorobantu VD, Le HM, Yue Y, Ames AD. 2019. *Episodic Learning with Control Lyapunov Functions for Uncertain Robotic Systems\**. In *2019 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, pp. 6878–6884. Piscataway, NJ: IEEE
120. Taylor A, Singletary A, Yue Y, Ames A. 2020. Learning for safety-critical control with control barrier functions. In *Proceedings of the 2nd Conference on Learning for Dynamics and Control*, ed. AM Bayen, A Jadbabaie, G Pappas, PA Parrilo, B Recht, C Tomlin, M Zeilinger, pp. 708–717, vol. 120 of *Proceedings of Machine Learning Research*. N.p.: PMLR
121. Ohnishi M, Wang L, Notomista G, Egerstedt M. 2019. Barrier-certified adaptive reinforcement learning with applications to brushbot navigation. *IEEE Transactions on Robotics* 35:1186–1205
122. Choi J, Castañeda F, Tomlin C, Sreenath K. 2020. *Reinforcement Learning for Safety-Critical Control under Model Uncertainty, using Control Lyapunov Functions and Control Barrier Functions*. In *Proceedings of Robotics: Science and Systems XVI*, pp. pap. 88. N.p.: Robot. Sci. Syst. Found.
123. Taylor AJ, Dorobantu VD, Krishnamoorthy M, Le HM, Yue Y, Ames AD. 2019. *A Control Lyapunov Perspective on Episodic Learning via Projection to State Stability*. In *2019 IEEE 58th Conference on Decision and Control (CDC)*, pp. 1448–1455. Piscataway, NJ: IEEE
124. Taylor AJ, Singletary A, Yue Y, Ames AD. 2020. A control barrier perspective on episodic learning via projection-to-state safety. [arXiv:2003.08028](https://arxiv.org/abs/2003.08028) [eess.SY]
125. Taylor AJ, Dorobantu VD, Dean S, Recht B, Yue Y, Ames AD. 2020. Towards robust data-driven control synthesis for nonlinear systems with actuation uncertainty. [arXiv:2011.10730](https://arxiv.org/abs/2011.10730) [eess.SY]
126. Taylor AJ, Ames AD. 2020. *Adaptive Safety with Control Barrier Functions*. In *2020 American Control Conference (ACC)*, pp. 1399–1405. Piscataway, NJ: IEEE
127. Lopez BT, Slotine JJE, How JP. 2021. Robust adaptive control barrier functions: An adaptive and data-driven approach to safety. *IEEE Control Systems Letters* 5:1031–1036
128. Cheng R, Orosz G, Murray RM, Burdick JW. 2019. *End-to-end safe reinforcement learning through barrier functions for safety-critical continuous control tasks*. In *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 33(1), pp. 3387–3395. Palo Alto, CA: AAAI Press
129. Fan DD, Nguyen J, Thakker R, Alatur N, Agha-mohammadi Aa, Theodorou EA. 2019. Bayesian learning-based adaptive control for safety critical systems. [arXiv:1910.02325](https://arxiv.org/abs/1910.02325) [eess.SY]
130. Wang L, Theodorou EA, Egerstedt M. 2018. *Safe learning of quadrotor dynamics using barrier certificates*. In *2018 IEEE International Conference on Robotics and Automation (ICRA)*, pp. 2460–2465. Piscataway, NJ: IEEE
131. Khojasteh MJ, Dhiman V, Franceschetti M, Atanasov N. 2020. Probabilistic safety constraints for learned high relative degree system dynamics. In *Proceedings of the 2nd Conference on Learning for Dynamics and Control*, ed. AM Bayen, A Jadbabaie, G Pappas, PA Parrilo, B Recht, C Tomlin, M Zeilinger, pp. 781–792, vol. 120 of *Proceedings of Machine Learning Research*. N.p.: PMLR
132. Mitchell I, Bayen A, Tomlin C. 2005. A time-dependent Hamilton-Jacobi formulation of reachable sets for continuous dynamic games. *IEEE Transactions on Automatic Control* 50:947–957
133. Fisac JF, Akametalu AK, Zeilinger MN, Kaynama S, Gillula J, Tomlin CJ. 2019. A general safety framework for learning-based control in uncertain robotic systems. *IEEE Transactions on Automatic Control* 64:2737–2752
134. Gillula JH, Tomlin CJ. 2012. *Guaranteed safe online learning via reachability: tracking a ground target using a quadrotor*. In *2012 IEEE International Conference on Robotics and Automation (ICRA)*, pp. 2723–2730. Piscataway, NJ: IEEE
135. Bajcsy A, Bansal S, Bronstein E, Tolani V, Tomlin CJ. 2019. *An Efficient Reachability-Based Framework for Provably Safe Autonomous Navigation in Unknown Environments*. In *2019*



- IEEE 58th Conference on Decision and Control (CDC)*, pp. 1758–1765. Piscataway, NJ: IEEE
136. Fisac JF, Lugovoy NF, Rubies-Royo V, Ghosh S, Tomlin CJ. 2019. *Bridging Hamilton-Jacobi safety analysis and reinforcement learning*. In *2019 International Conference on Robotics and Automation (ICRA)*, pp. 8550–8556. Piscataway, NJ: IEEE
  137. Choi JJ, Lee D, Sreenath K, Tomlin CJ, Herbert SL. 2021. Robust control barrier-value functions for safety-critical control. [arXiv:2104.02808](https://arxiv.org/abs/2104.02808) [eess.SY]
  138. Herbert S, Choi JJ, Sanjeev S, Gibson M, Sreenath K, Tomlin CJ. 2021. Scalable learning of safety guarantees for autonomous systems using Hamilton-Jacobi reachability. [arXiv:2101.05916](https://arxiv.org/abs/2101.05916) [cs.RO]
  139. Wabersich KP, Zeilinger MN. 2018. *Linear Model Predictive Safety Certification for Learning-Based Control*. In *2018 IEEE Conference on Decision and Control (CDC)*, pp. 7130–7135. Piscataway, NJ: IEEE
  140. Wabersich KP, Hewing L, Carron A, Zeilinger MN. 2019. Probabilistic model predictive safety certification for learning-based control. [arXiv:1906.10417](https://arxiv.org/abs/1906.10417) [eess.SY]
  141. Wabersich KP, Zeilinger MN. 2021. A predictive safety filter for learning-based control of constrained nonlinear dynamical systems. *Automatica* 129:109597
  142. Mannucci T, van Kampen E, de Visser C, Chu Q. 2018. Safe exploration algorithms for reinforcement learning controllers. *IEEE Transactions on Neural Networks and Learning Systems* 29:1069–1081
  143. Dean S, Taylor AJ, Cosner RK, Recht B, Ames AD. 2020. Guaranteeing safety of learned perception modules via measurement-robust control barrier functions. [arXiv:2010.16001](https://arxiv.org/abs/2010.16001) [eess.SY]
  144. Liu CK, Negrut D. 2021. The role of physics-based simulators in robotics. *Annual Review of Control, Robotics, and Autonomous Systems* 4:35–58
  145. Brockman G, Cheung V, Pettersson L, Schneider J, Schulman J, et al. 2016. OpenAI Gym. [arXiv:1606.01540](https://arxiv.org/abs/1606.01540) [cs.LG]
  146. Panerati J, Zheng H, Zhou S, Xu J, Prorok A, Schoellig AP. 2021. *Learning to Fly—a Gym Environment with PyBullet Physics for Reinforcement Learning of Multi-agent Quadcopter Control*. In *2021 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*. Piscataway, NJ: IEEE
  147. Coumans E, Bai Y. 2016–2021. PyBullet, a Python module for physics simulation for games, robotics and machine learning. <http://pybullet.org>
  148. Andersson JAE, Gillis J, Horn G, Rawlings JB, Diehl M. 2019. CasADi – A software framework for nonlinear optimization and optimal control. *Mathematical Programming Computation* 11:1–36
  149. Schulman J, Wolski F, Dhariwal P, Radford A, Klimov O. 2017. Proximal policy optimization algorithms. [arXiv:1707.06347](https://arxiv.org/abs/1707.06347) [cs.LG]
  150. Wieber PB, Tedrake R, Kuindersma S. 2016. Modeling and control of legged robots. In *Springer Handbook of Robotics*, ed. B Siciliano, O Khatib, pp. 1203–1234. Cham: Springer International Publishing
  151. McKinnon CD, Schoellig AP. 2018. *Experience-Based Model Selection to Enable Long-Term, Safe Control for Repetitive Tasks Under Changing Conditions*. In *2018 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, pp. 2977–2984. Piscataway, NJ: IEEE
  152. Chandak Y, Jordan S, Theodorou G, White M, Thomas PS. 2020. Towards safe policy improvement for non-stationary MDPs. In *Advances in Neural Information Processing Systems*, ed. H Larochelle, M Ranzato, R Hadsell, MF Balcan, H Lin, pp. 9156–9168, vol. 33. Red Hook, NY, USA: Curran Associates, Inc.
  153. Burgner-Kahrs J, Rucker DC, Choset H. 2015. Continuum robots for medical applications: A survey. *IEEE Transactions on Robotics* 31:1261–1280
  154. Mueller FL, Schoellig AP, D’Andrea R. 2012. *Iterative learning of feed-forward corrections for high-performance tracking*. In *2012 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, pp. 3276–3281. Piscataway, NJ: IEEE
  155. Dean S, Tu S, Matni N, Recht B. 2018. Safely learning to control the constrained linear quadratic regulator. *arXiv* :5582–5588
  156. McKinnon CD, Schoellig AP. 2019. *Learning Probabilistic Models for Safe Predictive Control in Unknown Environments*. In *2019 18th European Control Conference (ECC)*, pp. 2472–2479. Piscataway, NJ: IEEE

## A. SUMMARY OF REVIEWED LITERATURE

**Table 1** summarizes the safe learning control approaches reviewed in this paper. Learning approaches are classified as *Model-Based* if a dynamics model is used to produce control inputs or *Model-Free* if there is a direct mapping from states or measurements to control inputs. For model-based approaches, we categorize the structure of the *A Priori* and the *Learned* component of the model as Gaussian process (*GP*), neural network (*NN*), or *Other* for non-standard methods (e.g., set computation from data).

The *Safety Properties* category highlights which algorithms can handle state and/or input constraints (*Constraint Satisfaction*) and which ones provide stability guarantees (*Stability*). For both of these properties, we report whether they are enforced while the policy is being trained (*During Learning*) or by the resultant policy (*After Learning*). The *Robustness* category indicates whether the final policy is designed to be robust to *Input* disturbances; model *Parameter* uncertainty, and/or other *Dynamics* uncertainties broader than the previous two types of disturbances.

Finally, the last category reports the *Task(s)* to which each method was applied: *Stabilization* for stabilizing an equilibrium of the system; *Tracking* for the tracking of a given trajectory; *Navigation* for navigating or planning a sequence of actions to reach a given goal; *Locomotion* for the movement of legged robots like humanoids or quadrupeds; and *Manipulation* for pushing, reaching, or grasping operations using a robotic arm.

Table 1: Summary of the key properties of the works referenced in Sec. 3

Section	Learning			Safety Properties					Task	
	Model-based		Model-free	During Learning		After Learning				
	<i>A Priori</i>	Learned		Constraint Satisfaction	Stability	Constraint Satisfaction	Stability	Robustness w.r.t.		
Learning Uncertain Dynamics Sec. 3.1	Sec. 3.1.1	Linear (39, 40, 44, 45), Nonlinear (42, 43)	NN (39, 44, 45), GP (40, 42, 43)			(39, 40, 42, 43, 44, 45)		(39, 40, 42, 43, 44, 45)	Input (39, 40), Dynamics (42, 43, 44, 45)	Tracking (39, 40, 42, 43, 44, 45)
	Sec. 3.1.2	Linear (46, 47, 48), Nonlinear (49, 50)	GP (46, 48, 49, 50), Other (47)			(46, 48, 49, 50)		(46, 47, 48, 49, 50)	Dynamics (46, 47, 48, 49, 50)	Stabilization (46, 47, 48), Tracking (49, 50)
	Sec. 3.1.3	Linear (51, 52, 53, 54, 58, 59, 60, 61), Nonlinear (55, 56, 62, 63, 64, 65, 66, 67)	NN (60, 66), GP (61, 62, 63, 64, 65), Other (51, 52, 53, 54, 55, 56, 58, 59, 67)		Input & state (51, 52, 53, 54, 55, 56, 58, 59, 60, 61, 62, 63, 64, 65, 66)	(51, 52, 53, 54, 55, 56, 58, 59, 60, 61, 62, 63, 64, 65, 66)	Input & state (51, 52, 53, 54, 55, 56, 58, 59, 60, 61, 62, 63, 64, 65, 66, 67)	(51, 52, 53, 54, 55, 56, 58, 59, 60, 61, 62, 63, 64, 65, 66, 67)	Input (59), Dynamics (60, 61, 62, 63, 64, 65, 66, 67), Parameters (51, 52, 53, 54, 55, 56, 58, 59)	Stabilization (52, 53, 54, 55, 56, 58, 59, 60, 61, 64, 65), Tracking (51, 62, 63, 66, 67)
	Sec. 3.1.4	Nonlinear (68)	GP (68)			(68)		(68)		Stabilization (68)
Encouraging Safety in RL Sec. 3.2	Sec. 3.2.1	Nonlinear (69, 79, 82),	GP (78, 79, 80, 81)	NN (70, 75, 85, 86, 87), Tabular (73)	Input & state (69, 70, 73, 75, 78, 79, 80, 81, 82) Trajectory (87)		Input & state (69, 70, 73, 75, 78, 79, 80, 81, 82) Trajectory (85, 86, 87)			Stabilization (81), Tracking (79), Navigation (69, 70, 73, 82, 85, 86), Locomotion (85, 87), Manipulation (75, 85, 86, 87)
	Sec. 3.2.2	Nonlinear (89, 90)	NN (89, 90, 91, 92)	NN (93)			Input & state (89, 90, 91, 92)		Dynamics (93)	Stabilization (91), Navigation (89, 90, 91, 92), Locomotion (91, 93), Manipulation (92)
	Sec. 3.2.3	Tabular (97, 99)		NN (34, 95, 97, 98, 99)	Trajectory (34, 97, 98, 99)		Trajectory (34, 95, 97, 98, 99)			Navigation (34, 95, 97, 98, 99), Locomotion (34, 98, 99)
	Sec. 3.2.4		NN (106)	NN (36, 37, 38, 104, 105, 107, 108)	Input & state (106)	(106)	Input & state (106)	(106)	Input (37, 38), Dynamics (36, 104, 105, 106), Parameters (107, 108)	Stabilization (36, 104, 108), Tracking (106), Navigation (104, 105), Locomotion (36, 37, 38, 107), Manipulation (108)
Certifying Learning-based Control Sec. 3.3	Sec. 3.3.1	Linear (110), Nonlinear (109, 111, 113, 114)	NN (109, 110, 111, 113), Other (114)			(110)		(109, 110, 111, 113, 114)	Dynamics (109, 110, 111, 113, 114)	Stabilization (113, 114), Tracking (109, 110, 111, 114), Locomotion (114)
	Sec. 3.3.2	Linear (139, 140), Nonlinear (119, 120, 121, 122, 123, 124, 125, 126, 127, 128, 129, 130, 133, 134, 135, 141, 142)	GP (121, 128, 130, 131, 133, 140), NN (119, 120, 122, 123, 124, 128, 129, 135), Other (125, 126, 127, 134, 139, 140, 141)	NN (121, 122, 136)	Input & state (123, 124, 125, 126, 127, 128, 129, 130, 139, 140, 141, 142), State (133, 134, 135)	(123, 124, 125, 126, 127, 128, 129, 130, 133, 134, 135, 139, 140, 141)	Input & state (119, 120, 121, 122, 123, 124, 125, 126, 127, 128, 129, 130, 131, 139, 140, 141, 142), State (133, 134, 135, 136)	(119, 120, 121, 122, 123, 124, 125, 126, 127, 128, 129, 130, 131, 133, 134, 135, 136, 139, 140, 141)	Dynamics (119, 120, 121, 122, 123, 124, 125, 126, 128, 129, 130, 133, 134, 135, 136, 139, 140, 141)	Stabilization (128, 136, 139, 141, 142), Tracking (119, 120, 123, 124, 125, 126, 127, 129, 130, 131, 133, 134, 140), Navigation (121, 135), Locomotion (122)

## B. REVISION HISTORY

Send errata to [angela.schoellig@robotics.utias.utoronto.ca](mailto:angela.schoellig@robotics.utias.utoronto.ca).

Date	Change
December 6, 2021	Figure 6 update
August 13, 2021	Initial submission