

safe-control-gym: a Unified Benchmark Suite for Safe Learning-based Control and Reinforcement Learning in Robotics

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

Abstract—In recent years, both reinforcement learning and learning-based control—as well as the study of their *safety*, which is crucial for deployment in real-world robots—have gained significant traction. However, to adequately gauge the progress and applicability of new results, we need the tools to equitably compare the approaches proposed by the controls and reinforcement learning communities. Here, we propose a new open-source benchmark suite, called **safe-control-gym**, supporting both model-based and data-based control techniques. We provide implementations for three dynamic systems—the cart-pole, the 1D, and 2D quadrotor—and two control tasks—stabilization and trajectory tracking. We propose to extend OpenAI’s *Gym* API—the *de facto* standard in reinforcement learning research—with (i) the ability to specify (and query) symbolic dynamics and (ii) constraints, and (iii) (repeatedly) inject simulated disturbances in the control inputs, state measurements, and inertial properties. To demonstrate our proposal and in an attempt to bring research communities closer together, we show how to use **safe-control-gym** to quantitatively compare the control performance, data efficiency, and safety of multiple approaches from the fields of traditional control, learning-based control, and reinforcement learning.

CODE REPOSITORY

The open-source code repository of this project is at: <https://github.com/utiasDSL/safe-control-gym>.

I. INTRODUCTION

In a near future, robots will be the backbone of transportation, warehousing, and manufacturing. However, for ubiquitous robotics to materialize, we need to devise methods to develop robotic controllers faster and for increasingly complex systems—leveraging data and machine learning to scale up current design approaches. Top computing hardware and software companies (including Nvidia [1], Google [2], and intrinsic [3]) are currently working on fast physics-based simulation to complement real-world data and accelerate the progress of robot learning. At the same time, because safety is a crucial component of cyber-physical systems operating in the real world, safe learning-based control and safe reinforcement learning (RL) have become bustling areas of academic research over the past few years [4].

While the continuous influx of new contributions in safe robot learning reflects the relevance and progress of the field, we need to establish ways to fairly compare results yielded

The authors are with the Dynamic Systems Lab (DSL), Institute for Aerospace Studies, University of Toronto; the University of Toronto Robotics Institute; and affiliated with the Vector Institute for Artificial Intelligence in Toronto. E-mails: {firstname.lastname}@robotics.utias.utoronto.ca

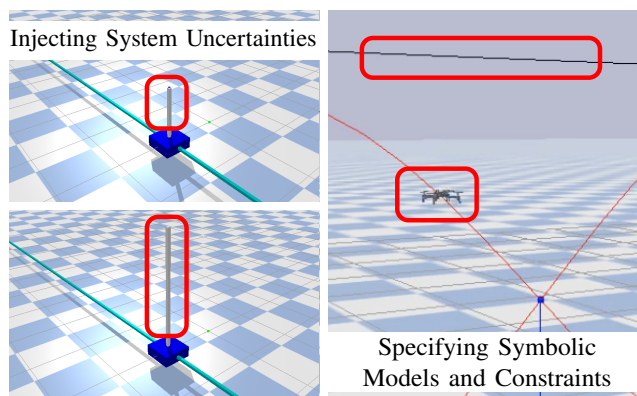


Fig. 1. **safe-control-gym** includes symbolic dynamics and constraints, and allows the injection of disturbances on inputs, states, and inertial properties to support the development of safe control algorithms for both stabilization and trajectory tracking tasks.

by controllers that leverage very different methodologies to understand their advantages and disadvantages. This is particularly important as different safe control approaches make different assumptions about the availability of prior knowledge of the system and of experimental data. Furthermore, we need quantitative benchmarks to assess these controllers’ safety and robustness.

Our work was motivated by the lack of open-source, reusable tools for the comparison of learning-based control research and RL (including methods that leverage prior knowledge) [4]. While we acknowledge the importance of eventual real-robot experiments, here, we focus on simulation as a way to lower the barrier of entry and appeal to a larger fraction of the control and RL communities. To develop safe learning-based robot control, we need a simulation environment that can (i) support model-based approaches, (ii) express safety constraints, and (iii) capture real-world non-idealities (such as uncertain physical properties and imperfect state estimation). Our ambition is for this software to bring closer, support, and speed up the work of control and RL researchers, allowing them to easily compare results. We strive for simple, modular, and reusable code which leverages two open-source tools popular with each of the two communities: PyBullet’s physics engine [5] and CasADi’s symbolic framework [6].

The contributions and features of the simulation environments in **safe-control-gym** (Figure 1) can be summarized as follows:

- we provide open-source simulation environments with a

TABLE I
FEATURE COMPARISON OF `safe-control-gym` AND OTHER SAFETY-FOCUSED REINFORCEMENT LEARNING ENVIRONMENTS

	Physics Engine	Rendering Engine	Robots	Tasks	Uncertain Conditions	Constraints	Disturbances	Gym API	Symbolic API
<code>safe-control-gym</code>	Bullet	TinyRenderer, OpenGL	Cart-pole, Quadrotor	Stabilization, Traj. Track.	Inertial Param., Initial State	State, Input	State, Input, Dynamics	Yes	Yes
<code>safety-gym</code> [7]	MuJoCo	OpenGL	Point, Car, Quadruped	Goal, Push, Button	Initial State	State	Adversaries	Yes	No
<code>realworldrl-suite</code> [8]	MuJoCo	OpenGL	Humanoid, Cart-pole, etc.	Swing-up, Walk, etc.	Inertial Param., Initial State	State	State, Input, Reward	No	No
<code>ai-safety-gridworlds</code> [9]	n/a	Terminal	n/a	Grid Navigation	Initial State, Reward	State	Dynamics, Adversaries	No	No

novel, augmented *Gym* API, including symbolic dynamics models and trajectory generation, designed to seamlessly interface with both RL and control approaches;

- `safe-control-gym` allows for constraint specification and disturbance injection onto a robot’s inputs, states, and inertial properties through a portable configuration system. This is crucial to simplify the development and comparison of safe learning-based control approaches;
- finally, our codebase includes open-source implementations of several baselines from traditional control, RL, and learning-based control, which we use in [4] and Section V to demonstrate how `safe-control-gym` supports quantitative comparisons across fields.

As pointed out in [4], [10], despite the undeniable similarities in their setup, there still exist terminology gaps and disconnects in how optimal control and reinforcement learning research address safe robot control. In [4], as we reviewed the last half-decade of safe robot control research, we observed significant differences in the use and reliance on prior models and assumptions. We also found a distinct lack of open-source simulations and control implementations (Figure 5 of [4]), which are essential for reproducibility and comparisons across fields and methodologies. With this work, we intend to make it easier for both RL and control researchers to (i) publish their results based on open-source simulations, (ii) compare against both RL and traditional control baselines, and (iii) quantify safety against shared, reusable sets of constraints and/or disturbances.

II. RELATED WORK

Simulation environments such as OpenAI’s `gym` [11] and DeepMind’s `dm_control` have been proposed as a way to standardize the development of RL algorithms and have been used, for example, in evaluation tool-kits for deep and model-based RL agents [12], [13]. However, recent work [14] has highlighted that, even using these tools, RL research is often difficult to reproduce as it might hinge on the choice of hyper-parameter values or random seeds.

In robotics research, physics-based simulators—such as Google’s Brax [2] and Nvidia’s Isaac Gym [1]—are an increasingly popular tool to effectively collect large data sets that approximate complex dynamical systems [15]. While MuJoCo has been the dominant force behind many of the physics-based RL environments, it is not yet a fully open-source project. In this work, we leverage the Python bindings of the open-source C++ Bullet Physics engine [5] instead,

which currently powers several re-implementations of the original MuJoCo tasks [16] as well as additional robotic simulations, including quadrotors [17] and quadrupeds.

Some aspects of safety have been investigated by existing simulation environments. DeepMind’s `ai-safety-gridworlds` [9], is a set of RL environments meant to assess the safety properties of intelligent agents by considering distributional shifts, robustness to adversaries, and safe exploration. However, it is not specific to real-world robotics systems, as these environments focus on grid worlds. OpenAI’s `safety-gym` [7] and Google’s `realworldrl-suite` [8] both augment typical RL environments with constraint evaluation. They include a handful of robotic platforms such as two-wheeled robots and quadrupeds. Similarly to our work, `realworldrl-suite` also includes perturbations in actions, observations, and physical quantities. However, unlike our work, the environments in [7], [8] lack the support for a symbolic framework to express *a priori* knowledge of a system’s dynamics or its constraints.

As our `safe-control-gym` includes a *Gym*-style quadrotor environment, it is important to clarify that we focus on safe, low-level control rather than vision-based applications (as AirSim [18] or Flightmare [19]) or multi-agent coordination (as `gym-pybullet-drones` [17]).

Our work advances the state-of-the-art in safety-focused simulation environments (summarized in Table I) by enabling the definition of: (i) symbolic models of the dynamics, cost, and constraints of an RL environment to support traditional control and model-based approaches; (ii) customizable, portable, and reusable constraints and physics disturbances to facilitate comparisons and enhance reproducibility; and (iii) several control and learning-based control baselines in addition to the typical RL baselines. These new features, paired with full retro-compatibility with *Gym*’s original API, make `safe-control-gym` an ideal playground to develop and compare RL and learning-based robot control.

III. ENVIRONMENTS

Our open-source suite `safe-control-gym` comprises three dynamical systems: the cart-pole, and the 1D and 2D quadrotors. We consider two control tasks: stabilization and trajectory tracking. The software package can be downloaded and installed as (see `README.md` for troubleshooting):

```
$ git clone -b submission \
> git@github.com:utiasDSL/safe-control-gym.git
```

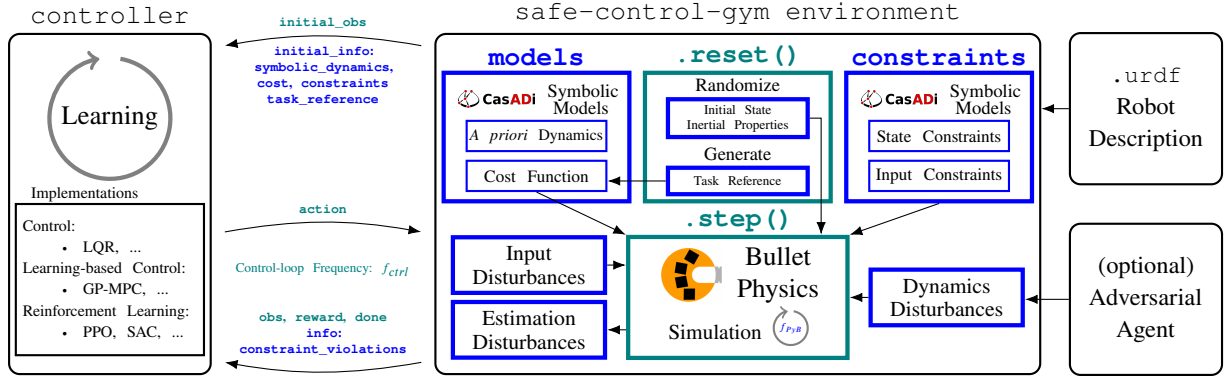


Fig. 2. Block diagram of `safe-control-gym`'s Python architecture, highlighting the backward compatibility with OpenAI's `Gym` API (teal components) and the proposed extensions (blue components) for the development of safe learning-based control and reinforcement learning.

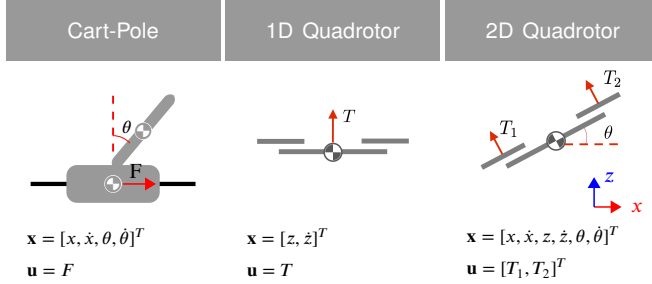


Fig. 3. Schematics, state and input vectors of the cart-pole, and the 1D and 2D quadrotor environments in `safe-control-gym`.

```
$ cd safe-control-gym/
$ pip3 install -e .
```

As advised in [20], “benchmark problems should be complex enough to highlight issues in controller design [...] but simple enough to provide easily understood comparisons.”. All three systems in `safe-control-gym` are unstable. We included the cart-pole, since it has been widely adopted to showcase traditional control as well as RL since the mid-80s [21]. The 1D quadrotor is a linear dynamic system, the 2D quadrotor is nonlinear, and the cart-pole is non-minimum phase. The 1D quadrotor is a simple double-integrator system, which can also be used for didactic purposes.

A. Cart-Pole System

An overview of the cart-pole system is given in Figure 3: a cart with mass m_c connects *via* a prismatic joint to a 1D track; a pole of mass m_p and length $2l$ is hinged to the cart. The state vector for the cart-pole is $\mathbf{x} = [x, \dot{x}, \theta, \dot{\theta}]^T$, where x is the horizontal position of the cart, \dot{x} is the velocity of the cart, θ is the angle of the pole with respect to vertical, and $\dot{\theta}$ is the angular velocity of the pole. The input to the system is a force $\mathbf{u} = F$ applied to the center of mass (COM) of the cart. In the frictionless case, the equations of motion for the cart-pole system are given in [22] as:

$$\ddot{x} = \frac{F + m_p l (\dot{\theta}^2 \sin \theta - \ddot{\theta} \cos \theta)}{m_c + m_p} \quad (1a)$$

$$\ddot{\theta} = \frac{g \sin \theta + \cos \theta \left(\frac{-F - m_p l \dot{\theta}^2 \sin \theta}{m_c + m_p} \right)}{l \left(\frac{4}{3} - \frac{m_p \cos^2 \theta}{m_c + m_p} \right)}, \quad (1b)$$

where g is the acceleration due to gravity.

B. 1D and 2D Quadrotor Systems

The second and third robotic systems in `safe-control-gym` are the 1D and the 2D quadrotor. These correspond to the cases in which the movement of a quadrotor is constrained to the 1D motion in the vertical z -direction and the 2D motion in the xz -plane, respectively (Figure 3). For a physical quadrotor, these motions can be achieved by controlling the four motor thrusts of the quadrotor to balance out the force and torque along the redundant dimensions.

In the 1D quadrotor case, the state of the system is $\mathbf{x} = [z, \dot{z}]^T$, where z and \dot{z} are the vertical position and velocity of the COM of the quadrotor. The input to the system is the overall thrust $\mathbf{u} = T$ generated by the motors of the quadrotor. The equation of motion for the 1D quadrotor system is:

$$\ddot{z} = T/m - g, \quad (2)$$

where m is the mass of the quadrotor and g is the acceleration due to gravity.

In the 2D quadrotor case, the state of the system is $\mathbf{x} = [x, \dot{x}, z, \dot{z}, \theta, \dot{\theta}]^T$, where (x, z) and (\dot{x}, \dot{z}) are the translation position and velocity of the COM of the quadrotor in the xz -plane, and θ and $\dot{\theta}$ are the pitch angle and the pitch angle rate, respectively. The input of the system are the thrusts $\mathbf{u} = [T_1, T_2]^T$ generated by two pairs of motors, one on each side of the body's y -axis. The equations of motion for the 2D quadrotor system are as follows:

$$\ddot{x} = \sin \theta (T_1 + T_2) / m \quad (3a)$$

$$\ddot{z} = \cos \theta (T_1 + T_2) / m - g \quad (3b)$$

$$\ddot{\theta} = (T_2 - T_1) d / I_{yy}, \quad (3c)$$

where m is the mass of the quadrotor, g is the acceleration due to gravity, $d = l/\sqrt{2}$ is the effective moment arm (with l being the arm length of the quadrotor, i.e., the distance from each motor pair to the COM), and I_{yy} is the moment of inertia about the y -axis.

C. Stabilization and Trajectory Tracking Tasks

For all three systems detailed in Sections III-A and III-B, our framework currently supports two control tasks: stabilization and trajectory tracking. For stabilization, `safe-control-gym` provides an equilibrium pair for the system, $\mathbf{x}^{\text{ref}}, \mathbf{u}^{\text{ref}}$. For trajectory tracking, we include a trajectory generation module capable of generating circular, sinusoidal, lemniscate, or square trajectories. The module returns references $\mathbf{x}_i^{\text{ref}}, \mathbf{u}_i^{\text{ref}} \forall i \in \{0, \dots, L\}$, where L is the number of control steps in an episode. To run a quadrotor example, tracking different trajectories with PID control, try:

```
$ cd safe-control-gym/examples/
$ python3 tracking.py --overrides ./tracking.yaml
```

For control-based approaches, on both stabilization and trajectory tracking tasks, `safe-control-gym` allows the computation of a quadratic cost J^Q of the form:

$$J^Q = \frac{1}{2} \sum_{i=0}^L (\mathbf{x}_i - \mathbf{x}_i^{\text{ref}})^T \mathbf{Q} (\mathbf{x}_i - \mathbf{x}_i^{\text{ref}}) + \frac{1}{2} \sum_{i=0}^{L-1} (\mathbf{u}_i - \mathbf{u}_i^{\text{ref}})^T \mathbf{R} (\mathbf{u}_i - \mathbf{u}_i^{\text{ref}}), \quad (4)$$

where \mathbf{Q}, \mathbf{R} are parameters of the cost function (note that, for the stabilization task, $\mathbf{x}_i^{\text{ref}}, \mathbf{u}_i^{\text{ref}}$ are not time-varying).

In RL, an agent’s performance is expressed by the total collected reward J^R and we use the negated quadratic cost to express it, that is $J^R = -J^Q$. The RL state \mathbf{x}_i is also augmented with the target position within the trajectory $\mathbf{x}_i^{\text{ref}}$. To enable further comparisons, we also implement the traditional reward function for cart-pole stabilization found in [11], [21]. This is an instantaneous reward of 1 for each discrete time step i in which the pole is upright, i.e., $|\theta| \leq \theta_{\max}$. Episodes are terminated on step D , when $|\theta|$ exceeds threshold θ_{\max} or $D = L$, thus $J^R = D$.

D. `safe-control-gym` Extended API

To provide native support to open-source RL libraries, `safe-control-gym` adopts OpenAI *Gym*’s interface. However, to the best of our knowledge, we are the first to extend this API with the ability to provide a learning agent with *a priori* knowledge of the dynamical system. This is of crucial importance to also support the development of, and comparison with, learning-based control approaches, which typically leverage insights about the physics of a robotic system. The proposed gym environment allows for this prior information to be integrated into the learning process. An overview of `safe-control-gym`’s modules is presented in Figure 2. Our benchmark suite can be used, for example, to answer the question of how data-efficiency, which is imperative for fast robot learning, is impacted by model-free RL approaches that do not exploit prior knowledge (see Section V). To run one of `safe-control-gym`’s environments, in headless mode, with terminal printouts from the original *Gym* API (in cyan) and our new API (in purple), try:

```
$ cd safe-control-gym/examples/
$ python3 verbose_api.py --system cartpole \
> --overrides verbose_api.yaml
```

The next four subsections detail the three core features of our simulations environments—*a priori* symbolic models, constraint specification, and disturbance injection—and the configuration system used for experiments’ reproducibility.

1) *Symbolic Models*: We use CasADi [6], an open-source symbolic framework for nonlinear optimization and algorithmic differentiation, to include symbolic models of (i) the systems’ *a priori* dynamics (see Sections III-A and III-B) as well as (ii) the quadratic cost function from Section III-C, and (iii) optional constraints (see Section III-D.2). As shown by the printouts of the snippet above, these models, together with the initial state of the system \mathbf{x}_0 and task references $\mathbf{x}^{\text{ref}}, \mathbf{u}^{\text{ref}}$, are returned by our API in a `reset_info` dictionary at each reset of an environment.

2) *Constraints*: The ability to specify, evaluate, and enforce one or more constraints $c^j \ j \in \{0, 1, 2, \dots\}$ on state \mathbf{x} and input \mathbf{u} :

$$c^j(\mathbf{x}_i, \mathbf{u}_i) \leq 0 \quad \forall i \in \{0, \dots, L\}, \quad (5)$$

is essential for safe robot control. While previous RL environments including state constraints exist [7], [8], our implementation is the first to also provide their symbolic representation and the ability to add bespoke ones while creating an environment (see Section III-D.4). Our current implementation includes default input and state constraints and supports user-specified ones in multiple forms (linear, bounded, quadratic) on either the system’s state, input, or both. Constraint evaluations are included in the `info` dictionary returned at each environment’s step (see Figure 2). With this constraint API, `safe-control-gym` can accommodate the different conventions and usage of constraints in safe control and safe RL, providing the necessary support to fairly evaluate and compare their safety properties.

3) *Disturbances*: In developing safe control approaches, we are often confronted with the fact that models like the ones in Sections III-A and III-B are not a complete or fully truthful representation of the system under test. `safe-control-gym` provides several ways to implement non-idealities that mimic real-life robots, covering many of the typical test scenarios in safe learning-based control as well as robust RL research, including:

- randomization (from a given probability distribution) of the initial state of the system, \mathbf{x}_0 ;
- randomization (from given probability distributions) of the inertial parameters, that is, m_c, m_p, l for the cart-pole and m, I_{yy} for the quadrotor;
- disturbances (in the form of white noise, step, or impulse) applied to the action input \mathbf{u} sent from the controller to the robot;
- disturbances (in the form of white noise, step, or impulse) applied to the observations of the state \mathbf{x} returned by an environment to the controller;
- external dynamics disturbances, additional forces applied to a robot using PyBullet APIs mimicking, for example, wind and other aerodynamic effects or set deterministically (from another agent outside the environment) to implement robust adversarial training schemes such as the one proposed in [23].

4) *Configuration System*: To facilitate the reproducibility and portability of experiments with identical environment setups, `safe-control-gym` includes a YAML configuration system that supports all of the features discussed in Sections III-C and III-D (see “Configuration” in `README.md`).

E. Computational Performance

The ability to collect large experimental datasets or generate representative simulated ones is one of the bottlenecks of learning-based robotics. With this in mind, we assessed the computational performance of `safe-control-gym` on a system with a 2.30GHz Quad-Core i7-1068NG7 CPU, 32GB 3733MHz LPDDR4X of memory, and running Python 3.8 under macOS 12. In the repository’s documentation, we summarize the simulation speed-ups (with respect to the wall-clock) obtained for the cart-pole and 2D quadrotor environments, in headless mode or using the GUI, with or without constraint evaluation, and different choices of control and physics integration frequencies. In headless mode, a single instance of `safe-control-gym` allows to collect data 10 to 20 times faster than in real life, with fine-grained physics stepped by PyBullet at 1000Hz.

IV. CONTROL ALGORITHMS

The codebase of `safe-control-gym` also comprises several implementations of core control approaches, from traditional and learning-based control, to safety-certified control and safe reinforcement learning [4].

A. Control and Safe Control Baselines

Our benchmark suite includes, as baselines, standard state-feedback controllers such as the linear quadratic regulator (LQR) and iterative LQR (iLQR) [24]. LQR assumes linear system dynamics (as in (2)) and quadratic cost (as in (4)). For nonlinear systems (e.g., the 2D quadrotor in (3) and cart-pole in (1)), LQR uses local linear approximations of the nonlinear dynamics. Unlike LQR, iLQR iteratively improves the performance by finding better local approximations of the cost function (4) and system dynamics using the state and input trajectories from the previous iteration. Our `safe-control-gym` environments expose the symbolic models of *a priori* dynamics and cost function, facilitating the computation of dynamics’ Jacobians and the Jacobians and Hessians of the cost function. We include LQR and iLQR to show how our benchmark supports model-based approaches and the symbolic expressions of the first-order and second-order terms included in each environment can be equivalently leveraged by other model-based methods.

We also implemented two predictive control baselines: linear model predictive control (LMPC) and nonlinear model predictive control (NMPC) [25]. At every control step, model predictive control (MPC) solves a constrained optimization problem to find a control input sequence, over a finite horizon, that minimizes the cost subject to the system’s predicted dynamics and possibly input and state constraints. Then, only the first optimal control input from the sequence is applied. While NMPC uses the nonlinear system model, LMPC uses the linearized approximation to predict the

evolution of the system, sacrificing prediction accuracy for computational efficiency. In our codebase, CasADi’s `opti` framework [6] is used to formulate the optimization problem. `safe-control-gym` provides all the system’s components required by MPC (*a priori* dynamics, constraints, cost function) as CasADi models, see Section III-D for details.

B. Reinforcement Learning Baselines

As `safe-control-gym` extends the original *Gym* API, any compatible RL algorithm can directly be applied to our environments. In our codebase, we include two of the most well-known RL baselines: Proximal Policy Optimization (PPO) [26] and Soft Actor-Critic (SAC) [27]. These are model-free approaches that map state measurements to control inputs without leveraging a dynamics model and using neural network (NN)-based policies. Both PPO and SAC have been shown to work on a wide range of simulated robotics tasks, some of which involve hybrid and nonlinear dynamics. We adapt their implementations from `stable-baselines3` [28] and OpenAI’s *Spinning Up*, with a few modifications to also support our suite’s configuration systems. PPO and SAC are not natively safety-aware approaches and do not guarantee constraint satisfaction nor robustness beyond the generalization properties of NNs.

C. Safe Learning-Based Control

Safe learning-based control approaches improve a robot’s performance using past data to improve the estimate of a system’s true dynamics. These approaches typically provide guarantees on stability and/or constraint satisfaction. One of these approaches, included in `safe-control-gym`, is GP-MPC [29]. This method models uncertain dynamics with a Gaussian process (GP), which is used to better predict the future evolution of the system as well as tighten constraints based on the confidence of the dynamics along the prediction horizon. GP-MPC has been demonstrated on off-road ground robots [29]. Our implementation leverages the NMPC controller, the environments’ symbolic *a priori* model, and uses `gpytorch` [30] for the GP modelling and optimization of the uncertain dynamics. GP-MPC can accommodate both environment and controller-specific constraints.

D. Safe and Robust Reinforcement Learning

Building upon the RL baselines, we implemented three safe RL approaches that address the problems of constraint satisfaction and robust generalization against disturbances from fixed distributions. The *safety layer*-based approach in [31] pre-trains NN models to approximate linearized state constraints. These learned constraints are then used to filter potentially unsafe inputs from an RL controller *via* least-squares projection (also see Section IV-E). Once pre-trained (without the need for a prior model), the safety layer is kept fixed during an agent’s learning. We add such a *safety layer* to PPO and apply it to our benchmark tasks. Robust RL aims to learn policies that generalize across systems or tasks. We adapt two methods based on adversarial learning: RARL [23] and RAP [32]. These approaches model disturbances as a learning adversary and train the policy against increasingly

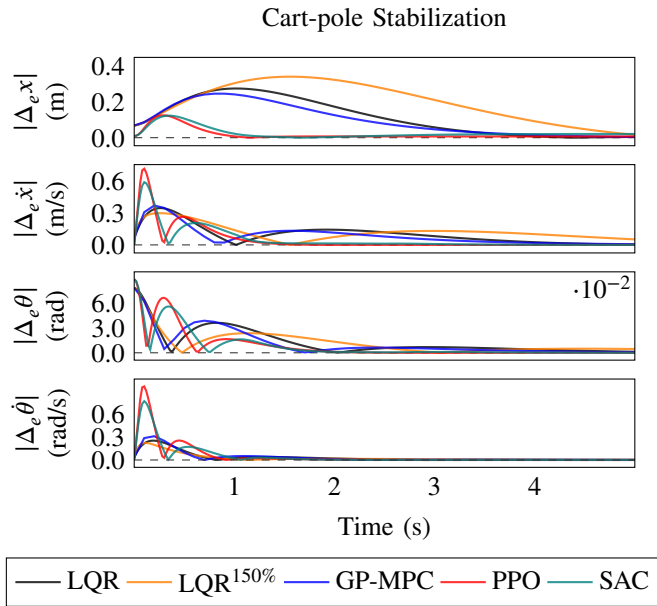


Fig. 4. Control performance (absolute error w.r.t. reference \mathbf{x}^{ref} , $|\Delta_e|$) on the cart-pole stabilization task for different controllers and RL agents.

stronger disturbances, which create harsher control scenarios. The controllers in [23], [32] demonstrated, in simulation, robustness against parameter mismatch and specific types of disturbances. These methods can be directly trained in *safe-control-gym*, by leveraging the disturbances API (see Section III-D.3).

E. Safety Certification of Learned Controllers

Learned controllers lacking formal guarantees themselves can be rendered safe by augmenting them with safety filters. These filters minimally modify unsafe control inputs, so that the applied control input maintains the system’s state within a safe set. A common safety filter that implements this idea is model predictive safety certification (MPSC), which solves a finite-horizon constrained optimization problem with a discrete-time predictive model to prevent a learning-based controller from violating constraints [33]. In [4], we presented an implementation of MPSC for PPO simultaneously leveraging the CasADi *a priori* dynamics and constraints and the *Gym* RL interface of *safe-control-gym*.

Control barrier functions (CBF) can act as safety filters for continuous-time nonlinear control-affine systems through a quadratic program (QP) with a constraint on the CBF’s time derivative [34]. In the case of model errors, the resulting errors in the CBF’s time derivative can be learned by a NN [35]. Learning-based CBF filters have been applied to safely control a segway [35] and a quadrotor [36]. Again, our CBF implementation relies on the *a priori* model and constraints exposed by *safe-control-gym*’s API. The CBF’s time derivative is also efficiently determined using CasADi. Constraints can be handled as long as the constraint set contains the safe set defined by the CBF.

V. RESULTS

In this section, we demonstrate how *safe-control-gym* can be used to compare approaches from all the families

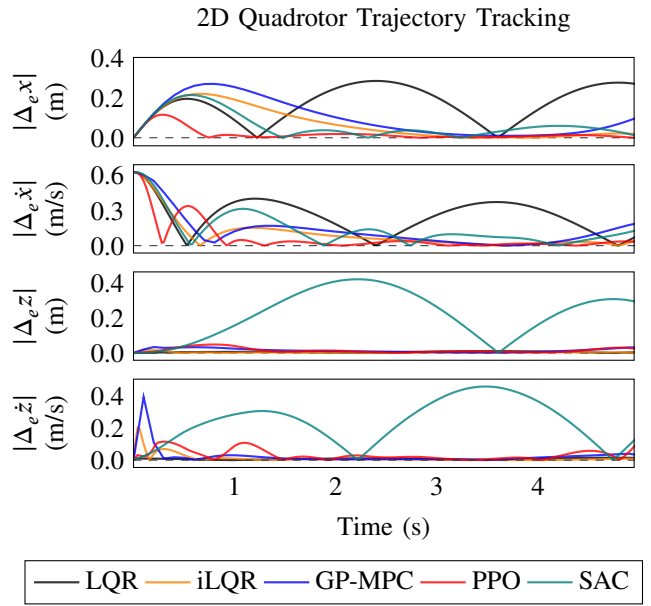


Fig. 5. Control performance (absolute error w.r.t. reference \mathbf{x}^{ref} , $|\Delta_e|$) on the 2D quadrotor tracking task for different controllers and RL agents.

of control algorithms discussed in Section IV, with respect to their control performance (Figures 4 and 5), learning efficiency (Figure 6), capability to satisfy constraints (Figure 7), as well as robustness against disturbances and parametric uncertainty (Figure 8). Our goal here is to highlight the potential for quantitative comparisons among controllers (although conditional to adequate and careful parameter tuning) as well as their qualitative features. Notably, *safe-control-gym* allows plotting RL and control results on a common set of axes with respect to the aforementioned performance and safety metrics, even when different controllers self-tune for different optimization objectives.

A. Control Performance

In Figures 4 and 5, we show that LQR, LQR^{150%}, GP-MPC (these latter two controllers using a prior model of the dynamics with parameters overestimated by 150%), PPO, and SAC are able to stabilize the cart-pole and track the quadrotor reference trajectory (a circle with a 1 meter radius). For the stabilization task, GP-MPC closely matches the closed-loop trajectory of the LQR with true parameters, albeit its *a priori* model was the same one given to LQR^{150%}. This shows how GP-MPC can overcome imperfect initial knowledge through learning. Both PPO and SAC yield substantially different closed-loop trajectories when compared to LQR and GP-MPC. Besides the fact that they have distinct learning procedures and parameterization, this is also due to the difference in their optimization objectives, and differences in prior information the algorithms have access to. Tracking the quadrotor sinusoidal trajectories (Figure 5) results in low-frequency oscillations for PPO and SAC, a consequence possibly from learning with stochastic data and the lack of explicit long-term planning. The LQR results in relatively large tracking errors due to the inaccuracy of the linearized dynamics model in approximating the nonlinear quadrotor system, especially for aggressive maneuvers.

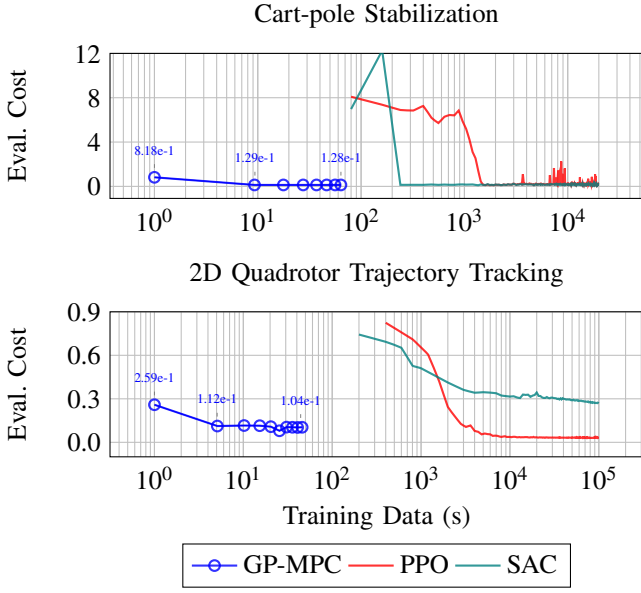


Fig. 6. Comparison of the learning performance (on a logarithmic x -axis) of learning-based control (GP-MPC) and two RL agents (PPO, SAC). Training data is expressed as seconds of simulation time used for collection. Evaluation cost is the RMSE with respect to the target state(s). In cart-pole stabilization, the RMSE uses the full state; in 2D quadrotor tracking, the RMSE only includes the x and z component of the state.

B. Learning Performance and Data Efficiency

Figure 6 shows how much data GP-MPC, PPO, and SAC require to achieve comparable trajectory root-mean-squared error (RMSE), a common metric for comparing controller performance. The values reported on the x -axis (training data) are measured in seconds of simulation time used for learning, derived as the number of elapsed environment simulation steps and PyBullet’s simulation frequency. This plot showcases the type of interdisciplinary comparisons uniquely enabled by `safe-control-gym`. In both plots, the untrained GP-MPC displays a performance that is only matched after around 10^3 seconds of simulated data by the RL approaches. GP-MPC converges to its optimal performance with roughly one tenth (or less) of the data. This highlights how learning-based control approaches are orders of magnitude more data-efficient than model-free RL. However, this is largely the result of knowing a reasonable *a priori* model. The evaluation costs of PPO and SAC exhibit large oscillations and learning instability in the early stage, not uncommon in deep RL [14]. Once converged SAC and PPO reach a performance comparable to GP-MPC in the stabilization task. PPO can match GP-MPC on the tracking task albeit using much more data. Although optimal hyper-parameter tuning is beyond the scope of this work, it should be noted that, for example, PPO’s learning stability is improved by a larger batch size, especially in the quadrotor tracking task.

C. Safety: Constraint Satisfaction

In Figure 7, we investigate the impact of learning and the amount of training data on the constraint violations of a learning-based controller or safe RL agent. The top plot summarizes the data efficiency of these approaches on the cart-pole stabilization task. Again, leveraging the

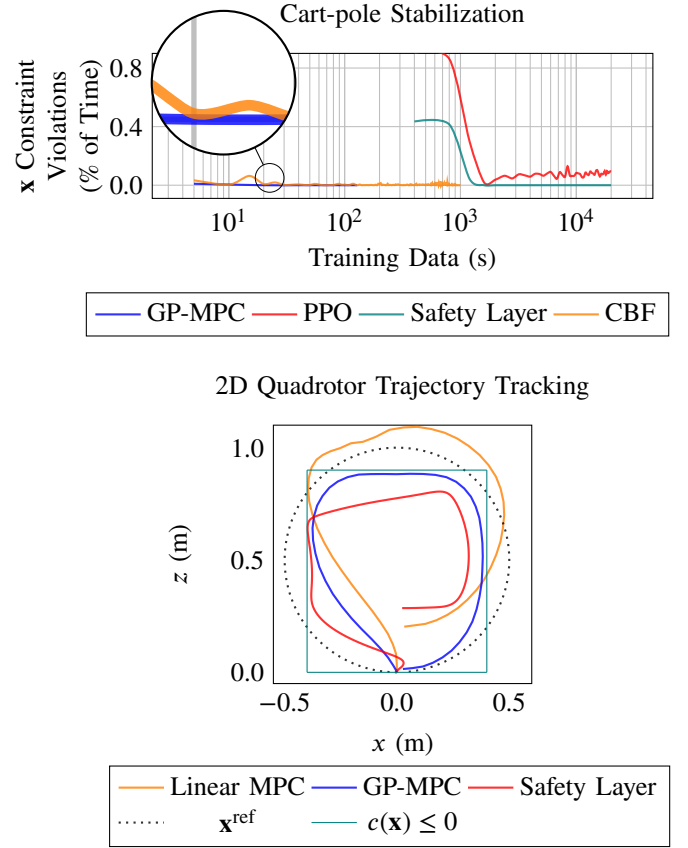


Fig. 7. In the top plot, fraction of time spent incurring a constraint violation by learning-based control (GP-MPC), vanilla RL (PPO), safety-augmented RL (Safety Layer), and safety-certified control (CBF); in the bottom plot, trajectories for an “impossible” tracking task (with constraints narrower than the reference) for traditional control (Linear MPC), learning-based control (GP-MPC), and safety-augmented RL (PPO with Safety Layer).

same overestimated prior model as in Section V-A, GP-MPC and the learning-based CBF require much fewer training examples to minimize the number of constraint violations than PPO with a *safety layer*. After training, GP-MPC, learning-based CBF, and safety layer PPO all achieve similar constraint satisfaction performance. Vanilla PPO also reduces the number of constraint violations but cannot match the performance of the GP-MPC and the learning-based CBF.

The bottom plot of Figure 7 shows reduced constraint violations for GP-MPC and PPO with a *safety layer* for the 2D quadrotor tracking task. Compared to a linear MPC with overestimated parameters, the GP-MPC meets the constraints, finding a compromise between performance and constraint satisfaction. PPO with a *safety layer*, on the other hand, struggles to balance between tracking the desired trajectory and fully guaranteeing constraint satisfaction.

D. Safety: Robustness

Leveraging `safe-control-gym`’s APIs for disturbances, Figure 8 shows the robustness of controllers and RL agents against parametric uncertainty (in the pole length) and white noise (on the input) for the cart-pole stabilization task. PPO^{DR} is trained with pole length randomization and improves the robustness of baseline PPO. RAP is trained against adversarial input disturbances and shows robust per-

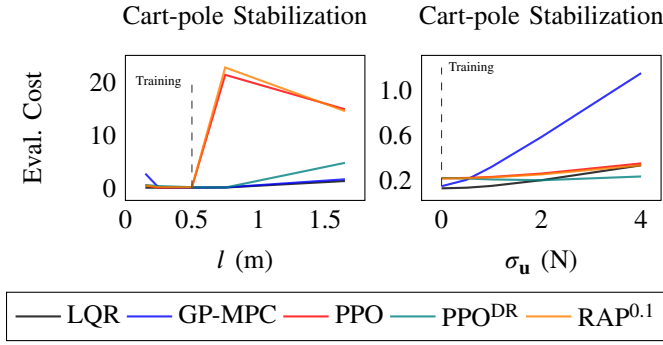


Fig. 8. Robustness of cart-pole stabilization policies learned by traditional (LQR) and learning-based control (GP-MPC) as well as vanilla (PPO) and robust RL (PPO^{DR}, RAP) with domain randomization against perturbations in the length of the pole l (left) and a white noise disturbance $\sim \mathcal{N}(0, \sigma_u^2)$ applied to action input \mathbf{u} (right).

formance to input noise, as expected from its adversarial training. However, RAP is less robust to different classes of disturbances, such as parameter mismatch. Model-based approaches, like LQR and GP-MPC, appear to be less affected by parameter uncertainty than model-free RL. This could be due to the fact that a prior dynamics model, albeit inexact, helps generalizing across similar systems. However, LQR and GP-MPC are no less resilient to input noise, when compared to PPO or RAP. GP-MPC, in particular, might be hindered by its lower control frequency (chosen to improve computational run-time), which results in accumulation of input noise that further degrades performance.

VI. CONCLUSIONS AND FUTURE WORK

In this letter, we introduced `safe-control-gym`, a benchmark suite of simulation environments to evaluate safe learning-based control. In `safe-control-gym`, we combine a physics engine-based simulation with the description of the available prior knowledge and safety constraints using a symbolic framework. By doing so, we allow the development and test of a wide range of approaches, from model-free RL to learning-based MPC. We hope that `safe-control-gym` will make it easier for researchers from the RL and control communities to compare their progress, especially for the quantification of safety and robustness in robotics applications.

REFERENCES

- [1] V. Makoviychuk, *et al.*, “Isaac gym: High performance GPU-based physics simulation for robot learning,” arXiv:2108.10470 [cs.RO], 2021.
- [2] C. D. Freeman, *et al.*, “Brax – a differentiable physics engine for large scale rigid body simulation,” arXiv:2106.13281 [cs.RO], 2021.
- [3] W. Tan-White, “Introducing intrinsic,” 2021. [Online]. Available: <https://blog.x.company/introducing-intrinsic-1cf35b87651>
- [4] L. Brunke, *et al.*, “Safe learning in robotics: From learning-based control to safe reinforcement learning,” *Annual Review of Control, Robotics, and Autonomous Systems*, vol. 5, no. 1, 2022. [Online]. Available: <https://doi.org/10.1146/annurev-control-042920-020211>
- [5] E. Coumans and Y. Bai, “PyBullet, a Python module for physics simulation for games, robotics and machine learning,” <http://pybullet.org>, 2016–2021.
- [6] J. A. E. Andersson, *et al.*, “CasADi – A software framework for non-linear optimization and optimal control,” *Mathematical Programming Computation*, vol. 11, no. 1, pp. 1–36, 2019.

- [7] A. Ray, *et al.*, “Benchmarking safe exploration in deep reinforcement learning,” <https://cdn.openai.com/safexp-short.pdf>, 2019.
- [8] G. Dulac-Arnold, *et al.*, “An empirical investigation of the challenges of real-world reinforcement learning,” arXiv:2003.11881 [cs.LG], 2021.
- [9] J. Leike, *et al.*, “AI safety gridworlds,” arXiv:1711.09883 [cs.LG], 2017.
- [10] B. Recht, “A tour of reinforcement learning: The view from continuous control,” *Annual Review of Control, Robotics, and Autonomous Systems*, vol. 2, no. 1, pp. 253–279, 2019.
- [11] G. Brockman, *et al.*, “OpenAI Gym,” arXiv:1606.01540 [cs.LG], 2016.
- [12] R. Julian, *et al.*, “Garage: A toolkit for reproducible reinforcement learning research,” <https://github.com/rllworkgroup/garage>, 2019.
- [13] T. Wang, *et al.*, “Benchmarking model-based reinforcement learning,” arXiv:1907.02057 [cs.LG], 2019.
- [14] P. Henderson, *et al.*, “Deep reinforcement learning that matters,” in *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 32(1). AAAI Press, 2018.
- [15] J. Collins, *et al.*, “A review of physics simulators for robotic applications,” *IEEE Access*, vol. 9, pp. 51 416–51 431, 2021.
- [16] B. Ellenberger, “PyBullet Gymperium,” <https://github.com/benelot/pybullet-gym>, 2018–2019.
- [17] J. Panerati, *et al.*, “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)*, 2021, pp. 7512–7519.
- [18] S. Shah, *et al.*, “AirSim: High-fidelity visual and physical simulation for autonomous vehicles,” in *Field and Service Robotics*. Springer Int’l Publishing, 2018, pp. 621–635.
- [19] Y. Song, *et al.*, “Flightmare: A flexible quadrotor simulator,” in *Proc. of the 4th Conference on Robot Learning*, 2020.
- [20] J. P. How, “Benchmarks [from the editor],” *IEEE Control Systems Magazine*, vol. 35, no. 1, pp. 6–7, 2015.
- [21] A. G. Barto, *et al.*, “Neuronlike adaptive elements that can solve difficult learning control problems,” *IEEE Transactions on Systems, Man, and Cybernetics*, vol. SMC-13, no. 5, pp. 834–846, 1983.
- [22] R. V. Florian, “Correct equations for the dynamics of the cart-pole system,” *Center for Cognitive and Neural Studies (Coneural)*, Romania, 2007.
- [23] L. Pinto, *et al.*, “Robust adversarial reinforcement learning,” in *Proceedings of the 34th International Conference on Machine Learning*, 2017, vol. 70, pp. 2817–2826.
- [24] J. Buchli, *et al.*, “Optimal and learning control for autonomous robots,” arXiv:1708.09342 [cs.SY], 2017.
- [25] J. B. Rawlings, *et al.*, *Model Predictive Control: Theory, Computation, and Design*. Nob Hill Publishing, 2020, vol. 2nd.
- [26] J. Schulman, *et al.*, “Proximal policy optimization algorithms,” arXiv:1707.06347 [cs.LG], 2017.
- [27] T. Haarnoja, *et al.*, “Soft actor-critic: Off-policy maximum entropy deep reinforcement learning with a stochastic actor,” in *Proceedings of the 35th International Conference on Machine Learning*, vol. 80, 2018, pp. 1861–1870.
- [28] A. Raffin, *et al.*, “Stable baselines3,” <https://github.com/DLR-RM/stable-baselines3>, 2019.
- [29] L. Hewing, *et al.*, “Cautious model predictive control using gaussian process regression,” *IEEE Transactions on Control Systems Technology*, vol. 28, no. 6, pp. 2736–2743, 2020.
- [30] J. R. Gardner, *et al.*, “GPyTorch: Blackbox matrix-matrix Gaussian process inference with GPU acceleration,” 2021.
- [31] G. Dalal, *et al.*, “Safe exploration in continuous action spaces,” arXiv:1801.08757 [cs.AI], 2018.
- [32] E. Vinitzky, *et al.*, “Robust reinforcement learning using adversarial populations,” arXiv:2008.01825 [cs.LG], 2020.
- [33] K. P. Wabersich and M. N. Zeilinger, “Linear model predictive safety certification for learning-based control,” in *2018 IEEE Conference on Decision and Control (CDC)*, 2018, pp. 7130–7135.
- [34] A. D. Ames, *et al.*, “Control barrier functions: Theory and applications,” in *2019 18th European Control Conference (ECC)*, 2019, pp. 3420–3431.
- [35] A. Taylor, *et al.*, “Learning for safety-critical control with control barrier functions,” in *Proceedings of the 2nd Conference on Learning for Dynamics and Control*, 2020, vol. 120, pp. 708–717.
- [36] L. Wang, *et al.*, “Safe learning of quadrotor dynamics using barrier certificates,” in *2018 IEEE International Conference on Robotics and Automation (ICRA)*, 2018, pp. 2460–2465.