

A Control Barrier Perspective on Episodic Learning via Projection-to-State Safety

Andrew J. Taylor, Andrew Singletary, Yisong Yue, Aaron D. Ames

Abstract—In this paper we seek to quantify the ability of learning to improve safety guarantees endowed by Control Barrier Functions (CBFs). In particular, we investigate how model uncertainty in the time derivative of a CBF can be reduced via learning, and how this leads to stronger statements on the safe behavior of a system. To this end, we build upon the idea of Input-to-State Safety (ISSf) to define Projection-to-State Safety (PSSf), which characterizes degradation in safety in terms of a projected disturbance. This enables the direct quantification of both how learning can improve safety guarantees, and how bounds on learning error translate to bounds on degradation in safety. We demonstrate that a practical episodic learning approach can use PSSf to reduce uncertainty and improve safety guarantees in simulation and experimentally.

I. INTRODUCTION

Ensuring safety is of significant importance in the design of many modern control systems, from autonomous driving to industrial robotics. In practice, the models used in the control design process are imperfect, with model uncertainty arising due to parametric error and unmodeled dynamics. This uncertainty can cause the controller to render the system unsafe. As such, it is necessary to quantify how the desired safety properties degrade with uncertainty.

Control Barrier Functions (CBFs) have become increasingly popular [15], [21], [2] as a tool for synthesizing controllers that provide safety via *set invariance* [6]. Safety guarantees endowed by a controller synthesized via CBFs rely on an accurate model of a system's dynamics, and may degrade in the presence of model uncertainty. The recently proposed definition of *Input-to-State Safety* (ISSf) provides a tool for quantifying the impact on safety guarantees of such uncertainty or disturbances in the dynamics [13] by describing changes in the set kept invariant.

Due to its flexibility, it is increasingly popular to incorporate learning into safe controller synthesis [22], [8], [16], [5], [9]. Many of these approaches seek to provide statistical guarantees on the safety via assumptions made on learning performance. In practice however, limitations on learning performance arise due to factors such as covariate shift [7], [14], limitations on model capacity, and optimization error. Thus, it is critical to understand the relationship between learning error and what safety guarantees can be ensured.

In this paper, we study how introducing learning models into safe controller synthesis done via CBFs can improve safety guarantees, and what safety guarantees can be made in

the presence of learning error. In particular, we consider the episodic learning approach proposed in [20], where learning is done directly on the time derivative of a CBF. We integrate this approach with Input-to-State Safety to not only highlight how learning can intuitively lead to improved safety guarantees, but also provide a direct relationship between learning error and the degradation of safety guarantees.

We make two main contributions in this paper. First, inspired by the idea of Projection-to-State Stability proposed in [19], we formulate general definitions of projections and projection compatible functions. Care must be taken to ensure these definitions preserve important topological properties for safety such as safe set membership. These definitions not only capture the definitions established in [19] as a special case, but allow us to define the notion of *Projection-to-State Safety* (PSSf), which is a variant of the Input-to-State Safety property. Like ISSf, PSSf provides a tool for characterizing the degradation of safety in the presence of disturbances. Unlike ISSf, PSSf considers disturbances in a projected environment, allowing stronger guarantees on safe behavior. Second, we demonstrate the utility of PSSf by characterizing how data-driven learning models can improve safety guarantees, and how learning error leads to degradation in safety guarantees.

Our paper is organized as follows. Section II provides a review of Control Barrier Functions and Input-to-State Safety. In Section III we define Projection-to-State Safety (PSSf) and discuss how PSSf enables quantifying degradation of safety in terms of a projected disturbance. Section IV defines a broad class of model uncertainty and explores how learning can be used to mitigate the impact of this uncertainty on safety. Lastly, in Section V we present both simulation and experimental results using PSSf to quantify the impact of learning error on safety guarantees for a Segway system.

II. PRELIMINARIES

This section provides a review of Control Barrier Functions (CBFs) and Input-to-State Safe Control Barrier Functions (ISSf-CBFs). These tools will be used in Section III to define the notion of Projection-to-State Safety.

Consider the nonlinear control affine system given by:

$$\dot{\mathbf{x}} = \mathbf{f}(\mathbf{x}) + \mathbf{g}(\mathbf{x})\mathbf{u}, \quad (1)$$

where $\mathbf{x} \in \mathbb{R}^n$, $\mathbf{u} \in \mathbb{R}^m$, and $\mathbf{f} : \mathbb{R}^n \rightarrow \mathbb{R}^n$ and $\mathbf{g} : \mathbb{R}^n \rightarrow \mathbb{R}^{n \times m}$ are locally Lipschitz continuous on \mathbb{R}^n . Given a Lipschitz continuous state-feedback controller $\mathbf{k} : \mathbb{R}^n \rightarrow \mathbb{R}^m$, the closed-loop system dynamics are:

$$\dot{\mathbf{x}} = \mathbf{f}_{\text{cl}}(\mathbf{x}) \triangleq \mathbf{f}(\mathbf{x}) + \mathbf{g}(\mathbf{x})\mathbf{k}(\mathbf{x}). \quad (2)$$

The assumption on local Lipschitz continuity of \mathbf{f} and \mathbf{k} implies that \mathbf{f}_{cl} is locally Lipschitz continuous. Thus for any

initial condition $\mathbf{x}_0 := \mathbf{x}(0) \in \mathbb{R}^n$ there exists a maximum time interval $I(\mathbf{x}_0) = [0, t_{\max})$ such that $\mathbf{x}(t)$ is the unique solution to (2) on $I(\mathbf{x}_0)$ [17]. In the case that \mathbf{f}_{cl} is forward complete, $t_{\max} = \infty$.

A continuous function $\alpha : [0, a) \rightarrow \mathbb{R}_+$, with $a > 0$, is said to belong to *class* \mathcal{K} ($\alpha \in \mathcal{K}$) if $\alpha(0) = 0$ and α is strictly monotonically increasing. If $a = \infty$ and $\lim_{r \rightarrow \infty} \alpha(r) = \infty$, then α is said to belong to *class* \mathcal{K}_∞ ($\alpha \in \mathcal{K}_\infty$). A continuous function $\alpha : (-b, a) \rightarrow \mathbb{R}$, with $a, b > 0$, is said to belong to *extended class* \mathcal{K} ($\alpha \in \mathcal{K}_e$) if $\alpha(0) = 0$ and α is strictly monotonically increasing. If $a, b = \infty$, $\lim_{r \rightarrow \infty} \alpha(r) = \infty$, and $\lim_{r \rightarrow -\infty} \alpha(r) = -\infty$, then α is said to belong to *extended class* \mathcal{K}_∞ ($\alpha \in \mathcal{K}_{\infty,e}$).

The notion of safety that we consider is formalized by specifying a *safe set* in the state space that the system must remain in to be considered safe. In particular, consider a set $\mathcal{C} \subset \mathbb{R}^n$ defined as the 0-superlevel set of a continuously differentiable function $h : \mathbb{R}^n \rightarrow \mathbb{R}$, yielding:

$$\mathcal{C} \triangleq \{\mathbf{x} \in \mathbb{R}^n : h(\mathbf{x}) \geq 0\}, \quad (3)$$

$$\partial\mathcal{C} \triangleq \{\mathbf{x} \in \mathbb{R}^n : h(\mathbf{x}) = 0\}, \quad (4)$$

$$\text{Int}(\mathcal{C}) \triangleq \{\mathbf{x} \in \mathbb{R}^n : h(\mathbf{x}) > 0\}. \quad (5)$$

We assume that \mathcal{C} is nonempty and has no isolated points, that is, $\text{Int}(\mathcal{C}) \neq \emptyset$ and $\overline{\text{Int}(\mathcal{C})} = \mathcal{C}$. We refer to \mathcal{C} as the *safe set*. This construction motivates the following definitions of forward invariant and safety:

Definition 1 (Forward Invariant & Safety). A set $\mathcal{C} \subset \mathbb{R}^n$ is *forward invariant* if for every $\mathbf{x}_0 \in \mathcal{C}$, the solution $\mathbf{x}(t)$ to (2) satisfies $\mathbf{x}(t) \in \mathcal{C}$ for all $t \in I(\mathbf{x}_0)$. The system (2) is *safe* on the set \mathcal{C} if the set \mathcal{C} is forward invariant.

Certifying the safety of the closed-loop system (2) with respect to a set \mathcal{C} may be impossible if the controller \mathbf{k} was not chosen to enforce the safety of \mathcal{C} . Control Barrier Functions can serve as a synthesis tool for attaining the forward invariance, and thus the safety of a set:

Definition 2 (Control Barrier Function (CBF), [4]). Let $\mathcal{C} \subset \mathbb{R}^n$ be the 0-superlevel set of a continuously differentiable function $h : \mathbb{R}^n \rightarrow \mathbb{R}$ with 0 a regular value. The function h is a *Control Barrier Function (CBF)* for (1) on \mathcal{C} if there exists $\alpha \in \mathcal{K}_{\infty,e}$ such that for all $\mathbf{x} \in \mathbb{R}^n$:

$$\sup_{\mathbf{u} \in \mathbb{R}^m} \dot{h}(\mathbf{x}, \mathbf{u}) \triangleq \frac{\partial h}{\partial \mathbf{x}}(\mathbf{x})(\mathbf{f}(\mathbf{x}) + \mathbf{g}(\mathbf{x})\mathbf{u}) \geq -\alpha(h(\mathbf{x})). \quad (6)$$

We note that this definition can be relaxed such that the inequality only holds for all $\mathbf{x} \in E$ where E is an open set satisfying $\mathcal{C} \subset E \subset \mathbb{R}^n$. Given a CBF h for (1) and a corresponding $\alpha \in \mathcal{K}_{\infty,e}$, we can consider the point-wise set of all control values that satisfy (6):

$$K_{\text{cbf}}(\mathbf{x}) \triangleq \left\{ \mathbf{u} \in \mathbb{R}^m \mid \dot{h}(\mathbf{x}, \mathbf{u}) \geq -\alpha(h(\mathbf{x})) \right\}.$$

One of the main results in [1], [23] relates controllers taking values in $K_{\text{cbf}}(\mathbf{x})$ to the safety of (1) on \mathcal{C} :

Theorem 1. *Given a set $\mathcal{C} \subset \mathbb{R}^n$ defined as the 0-superlevel set of a continuously differentiable function $h : \mathbb{R}^n \rightarrow \mathbb{R}$, if h is a CBF for (1) on \mathcal{C} , then any Lipschitz continuous*

controller $\mathbf{k} : \mathbb{R}^n \rightarrow \mathbb{R}^m$, such that $\mathbf{k}(\mathbf{x}) \in K_{\text{cbf}}(\mathbf{x})$ for all $\mathbf{x} \in \mathbb{R}^n$, renders the system (1) safe with respect to the set \mathcal{C} .

To accommodate disturbances or model uncertainties, we consider a disturbance space $\mathcal{D} \in \mathbb{R}^n$ and a disturbed system:

$$\dot{\mathbf{x}} = \mathbf{f}(\mathbf{x}) + \mathbf{g}(\mathbf{x})\mathbf{u} + \mathbf{d}. \quad (7)$$

with $\mathbf{d} \in \mathcal{D}$. The disturbance may be time-varying, state and/or input dependent. We will assume that when viewing \mathbf{d} as a signal, $\mathbf{d}(t)$, it is essentially bounded in time, and define $\|\mathbf{d}\|_\infty \triangleq \text{ess sup}_{t \geq 0} \|\mathbf{d}(t)\|$. Under a Lipschitz continuous state-feedback controller \mathbf{k} , the closed-loop dynamics are then given by:

$$\dot{\mathbf{x}} = \mathbf{f}_{\text{cl}}(\mathbf{x}, \mathbf{d}) \triangleq \mathbf{f}(\mathbf{x}) + \mathbf{g}(\mathbf{x})\mathbf{k}(\mathbf{x}) + \mathbf{d}. \quad (8)$$

In the presence of disturbances, a controller \mathbf{k} synthesized to render the set \mathcal{C} safe for the undisturbed dynamics (2) may fail to render \mathcal{C} safe for the disturbed dynamics (8). To quantify how safety degrades, we consider the notion of *input-to-state safety* [13]:

Definition 3 (Input-to-State Safety (ISSf)). The closed-loop system (8) is *input-to-state safe (ISSf)* on a set $\mathcal{C} \subset \mathbb{R}^n$ with respect to disturbances \mathbf{d} if there exists $\bar{d} > 0$ and $\gamma \in \mathcal{K}_\infty$ such that the set $\mathcal{C}_{\mathbf{d}} \supset \mathcal{C}$ defined as:

$$\mathcal{C}_{\mathbf{d}} \triangleq \{\mathbf{x} \in \mathbb{R}^n : h(\mathbf{x}) + \gamma(\|\mathbf{d}\|_\infty) \geq 0\}, \quad (9)$$

$$\partial\mathcal{C}_{\mathbf{d}} \triangleq \{\mathbf{x} \in \mathbb{R}^n : h(\mathbf{x}) + \gamma(\|\mathbf{d}\|_\infty) = 0\}, \quad (10)$$

$$\text{Int}(\mathcal{C}_{\mathbf{d}}) \triangleq \{\mathbf{x} \in \mathbb{R}^n : h(\mathbf{x}) + \gamma(\|\mathbf{d}\|_\infty) > 0\}, \quad (11)$$

is forward invariant for all \mathbf{d} satisfying $\|\mathbf{d}\|_\infty \leq \bar{d}$.

We refer to \mathcal{C} as an *input-to-state safe set (ISSf set)* if such a set $\mathcal{C}_{\mathbf{d}}$ exists. This definition implies that though the set \mathcal{C} may not be safe, a larger set $\mathcal{C}_{\mathbf{d}}$, depending on \mathbf{d} , is safe. If $\mathbf{d} \equiv \mathbf{0}$, we recover that the set \mathcal{C} is safe. \mathcal{C} can be certified as an ISSf set for the closed-loop system (8) with the following definition:

Definition 4 (Input-to-State Safe Barrier Function (ISSf-BF)). Let $\mathcal{C} \subset \mathbb{R}^n$ be the 0-superlevel set of a continuously differentiable function $h : \mathbb{R}^n \rightarrow \mathbb{R}$ with 0 a regular value. The function $h : \mathbb{R}^n \rightarrow \mathbb{R}$ is an *Input-to-State Safe Barrier Function (ISSf-BF)* for (8) on \mathcal{C} if there exists $\bar{d} > 0$, $\alpha \in \mathcal{K}_{\infty,e}$, and $\iota \in \mathcal{K}_\infty$ such that:

$$\frac{\partial h}{\partial \mathbf{x}}(\mathbf{x})(\mathbf{f}(\mathbf{x}) + \mathbf{g}(\mathbf{x})\mathbf{k}(\mathbf{x}) + \mathbf{d}) \geq -\alpha(h(\mathbf{x})) - \iota(\|\mathbf{d}\|), \quad (12)$$

for all $\mathbf{x} \in \mathbb{R}^n$ and $\mathbf{d} \in \mathbb{R}^n$ such that $\|\mathbf{d}\| \leq \bar{d}$.

As shown in [13], the existence of an ISSf-BF for (8) on \mathcal{C} implies \mathcal{C} is an ISSf set. Similarly to the undisturbed case, we can introduce the notion of a Control Barrier Function for synthesizing controllers that ensure input-to-state safety:

Definition 5 (ISSf Control Barrier Function (ISSf-CBF)). Let $\mathcal{C} \subset \mathbb{R}^n$ be the 0-superlevel set of a continuously differentiable function $h : \mathbb{R}^n \rightarrow \mathbb{R}$ with 0 a regular value. The function h is an *Input-to-State Safe Control Barrier Function (ISSf-CBF)* for (7) on \mathcal{C} if there exists $\bar{d} > 0$, $\alpha \in \mathcal{K}_{\infty,e}$, and $\iota \in \mathcal{K}_\infty$

such that:

$$\sup_{\mathbf{u} \in \mathbb{R}^m} \dot{h}(\mathbf{x}, \mathbf{u}, \mathbf{d}) \triangleq \frac{\partial h}{\partial \mathbf{x}}(\mathbf{x})(\mathbf{f}(\mathbf{x}) + \mathbf{g}(\mathbf{x})\mathbf{u} + \mathbf{d}) \geq -\alpha(h(\mathbf{x})) - \iota(\|\mathbf{d}\|), \quad (13)$$

for all $\mathbf{x} \in \mathbb{R}^n$ and $\mathbf{d} \in \mathbb{R}^n$ satisfying $\|\mathbf{d}\| \leq \bar{d}$.

We note that this definition is a more general definition of an ISSf-CBF compared to [13], where disturbances enter the system with the inputs. We define the pointwise set:

$$K_{\text{issf}}(\mathbf{x}) \triangleq \left\{ \mathbf{u} \in \mathbb{R}^m \mid \dot{h}(\mathbf{x}, \mathbf{u}, \mathbf{d}) \geq -\alpha(h(\mathbf{x})) - \iota(\|\mathbf{d}\|) \right\},$$

noting that for a fixed input the inequality must hold for all $\mathbf{d} \in \mathbb{R}^n$ satisfying $\|\mathbf{d}\| \leq \bar{d}$. Given this result, we have the following theorem:

Theorem 2. *Given a set $\mathcal{C} \subset \mathbb{R}^n$ defined as the 0-superlevel set of a continuously differentiable function $h : \mathbb{R}^n \rightarrow \mathbb{R}$, if h is an ISSf-CBF for (7) on \mathcal{C} , then any Lipschitz continuous controller $\mathbf{k} : \mathbb{R}^n \rightarrow \mathbb{R}^m$, such that $\mathbf{k}(\mathbf{x}) \in K_{\text{issf}}(\mathbf{x})$ for all $\mathbf{x} \in \mathbb{R}^n$, renders the set \mathcal{C} ISSf for (8).*

This theorem follows from the fact that under the controller \mathbf{k} , h serves an ISSf-BF for (8) on \mathcal{C} .

III. PROJECTION-TO-STATE SAFETY

Input-to-State Safety describes how the safe set \mathcal{C} changes in terms of the disturbance as it appears in the state dynamics (see Definition 3 in Section II). This description does not easily permit analysis of how safety degrades when the disturbance is more easily characterized by its impact in a Barrier Function derivative. This limitation motivates Projection-to-State Safety (PSSf), which enables a characterization of safety in terms of a projected disturbance.

We refer to a continuously differentiable function $\mathbf{\Pi} : \mathbb{R}^n \rightarrow \mathbb{R}^k$ as a *projection*, and denote $\mathbf{y} = \mathbf{\Pi}(\mathbf{x})$. Considering the system governed by (7), the associated projected system is governed by the dynamics:

$$\dot{\mathbf{y}} = \mathbf{D}_{\mathbf{\Pi}}(\mathbf{x})(\mathbf{f}(\mathbf{x}) + \mathbf{g}(\mathbf{x})\mathbf{u}) + \mathbf{D}_{\mathbf{\Pi}}(\mathbf{x})\mathbf{d}, \quad (14)$$

where $\mathbf{D}_{\mathbf{\Pi}} : \mathbb{R}^n \rightarrow \mathbb{R}^{k \times n}$ denotes the Jacobian of $\mathbf{\Pi}$. As will be seen when quantifying the impact of model uncertainty and learning error in Section IV, if the disturbance can be partially characterized in terms of the state and input, we may rewrite the projected dynamics as:

$$\dot{\mathbf{y}} = \mathbf{f}_{\mathbf{y}}(\mathbf{x}) + \mathbf{g}_{\mathbf{y}}(\mathbf{x})\mathbf{u} + \boldsymbol{\delta}, \quad (15)$$

where $\mathbf{f}_{\mathbf{y}} : \mathbb{R}^n \rightarrow \mathbb{R}^k$ and $\mathbf{g}_{\mathbf{y}} : \mathbb{R}^n \rightarrow \mathbb{R}^{k \times m}$ are Lipschitz continuous on \mathbb{R}^n , and $\boldsymbol{\delta} \in \mathbb{R}^k$ is referred to as the *projected disturbance*. We note it is not explicitly necessary that the relationships $\mathbf{f}_{\mathbf{y}}(\mathbf{x}) = \mathbf{D}_{\mathbf{\Pi}}(\mathbf{x})\mathbf{f}(\mathbf{x})$, $\mathbf{g}_{\mathbf{y}}(\mathbf{x}) = \mathbf{D}_{\mathbf{\Pi}}(\mathbf{x})\mathbf{g}(\mathbf{x})$, and $\boldsymbol{\delta} = \mathbf{D}_{\mathbf{\Pi}}(\mathbf{x})\mathbf{d}$ hold, but are one possible relationship between the terms in (14) and (15). For the following results, we will assume that $\boldsymbol{\delta}$ is essentially bounded in time and define $\|\boldsymbol{\delta}\|_{\infty} \triangleq \text{ess sup}_{t \geq 0} \|\boldsymbol{\delta}(t)\|$. We are interested in relating behaviors of the projected system to the original system, motivating the following definition:

Definition 6 (Projection-to-State Safety). The closed-loop system (8) is *projection-to-state safe* (PSSf) on \mathcal{C} with respect to the projection $\mathbf{\Pi}$ and projected disturbances $\boldsymbol{\delta}$ if there exists $\bar{\delta} > 0$ and $\gamma \in \mathcal{K}_{\infty}$ such that the set $\mathcal{C}_{\boldsymbol{\delta}} \supset \mathcal{C}$,

$$\mathcal{C}_{\boldsymbol{\delta}} \triangleq \{\mathbf{x} \in \mathbb{R}^n : h(\mathbf{x}) + \gamma(\|\boldsymbol{\delta}\|_{\infty}) \geq 0\}, \quad (16)$$

$$\partial \mathcal{C}_{\boldsymbol{\delta}} \triangleq \{\mathbf{x} \in \mathbb{R}^n : h(\mathbf{x}) + \gamma(\|\boldsymbol{\delta}\|_{\infty}) = 0\}, \quad (17)$$

$$\text{Int}(\mathcal{C}_{\boldsymbol{\delta}}) \triangleq \{\mathbf{x} \in \mathbb{R}^n : h(\mathbf{x}) + \gamma(\|\boldsymbol{\delta}\|_{\infty}) > 0\}, \quad (18)$$

is forward invariant for all $\boldsymbol{\delta}$ satisfying $\|\boldsymbol{\delta}\|_{\infty} \leq \bar{\delta}$.

In contrast to the definition of ISSf which enlarges the safe set in terms of the disturbance \mathbf{d} , PSSf quantifies how the safe set enlarges in terms of the projected disturbance $\boldsymbol{\delta}$. To utilize safety guarantees implied by ISSf-CBFs for analyzing PSSf behavior, we require the following definition:

Definition 7 (Compatible Projection). A function $h_{\mathbf{\Pi}} : \mathbb{R}^k \rightarrow \mathbb{R}$ is said to be a *compatible projection* for the function $h : \mathbb{R}^n \rightarrow \mathbb{R}$ with respect to the projection $\mathbf{\Pi} : \mathbb{R}^n \rightarrow \mathbb{R}^k$ if there exists $\underline{\sigma}, \bar{\sigma} \in \mathcal{K}_{\infty, e}$ such that for all $\mathbf{x} \in \mathbb{R}^n$:

$$\underline{\sigma}(h(\mathbf{x})) \leq h_{\mathbf{\Pi}}(\mathbf{\Pi}(\mathbf{x})) \leq \bar{\sigma}(h(\mathbf{x})). \quad (19)$$

Remark 1. If h and $h_{\mathbf{\Pi}}$ are norms on \mathbb{R}^n and \mathbb{R}^k , respectively, then $\mathbf{\Pi}$ reduces to a *dynamic projection* as introduced in [19]. Whereas dynamic projections preserve the topological notion of a point between the state and projected spaces, compatible projections can preserve more interesting topological structures such as sets.

Remark 2. The definition of a compatible projection can be abstractly viewed through the lens of category theory, mirroring the idea that one proves a property by mapping a system to the ‘‘simplist’’ type of system that has that property [3]. For safety, these are dynamical systems defined on the entire real line, with the safe set being the positive reals. Thus $h_{\mathbf{\Pi}}$ is a compatible projection if the following diagram:

$$\begin{array}{ccc} \mathbb{R}^n & \xrightarrow{h} & \mathbb{R} \\ \mathbf{\Pi} \downarrow & \nearrow h_{\mathbf{\Pi}} & \\ \mathbb{R}^k & & \end{array}$$

commutes up to class \mathcal{K} functions, i.e., (19) being satisfied.

In the context of safety, if a set $\mathcal{C} \subset \mathbb{R}^n$ is defined via a continuously differentiable function h as in (3)-(5), a compatible projection $h_{\mathbf{\Pi}}$ for the function h with respect to $\mathbf{\Pi}$ defines a corresponding set $\mathcal{C}_{\mathbf{\Pi}} \subset \mathbb{R}^k$:

$$\mathcal{C}_{\mathbf{\Pi}} \triangleq \{\mathbf{y} \in \mathbb{R}^k : h_{\mathbf{\Pi}}(\mathbf{y}) \geq 0\}, \quad (20)$$

$$\partial \mathcal{C}_{\mathbf{\Pi}} \triangleq \{\mathbf{y} \in \mathbb{R}^k : h_{\mathbf{\Pi}}(\mathbf{y}) = 0\}, \quad (21)$$

$$\text{Int}(\mathcal{C}_{\mathbf{\Pi}}) \triangleq \{\mathbf{y} \in \mathbb{R}^k : h_{\mathbf{\Pi}}(\mathbf{y}) > 0\}. \quad (22)$$

The inequalities in (19) preserve the notion of what states are considered safe between the state space and projected space, such that $\mathbf{x} \in \mathcal{C} \implies \mathbf{\Pi}(\mathbf{x}) \in \mathcal{C}_{\mathbf{\Pi}}$. The preceding implication is also true of the boundaries and interiors of the two sets. The following theorem allows us to extend ISSf properties of the projected system on $\mathcal{C}_{\mathbf{\Pi}}$ to PSSf properties of the original system on \mathcal{C} .

Theorem 3. Let $\mathcal{C} \subset \mathbb{R}^n$ be the 0-superlevel set of a continuously differentiable function $h : \mathbb{R}^n \rightarrow \mathbb{R}$ with 0 a regular value. The disturbed system (7) can be rendered PSSf on \mathcal{C} with respect to the projection Π and projected disturbances δ if there exists a compatible projection h_{Π} for h with respect to Π and Lipschitz continuous controller $\mathbf{k} : \mathbb{R}^n \rightarrow \mathbb{R}^m$ such that h_{Π} is an ISSf-CBF for the projected dynamics (15) on \mathcal{C}_{Π} and $\mathbf{k}(\mathbf{x}) \in K_{\text{issf}}(\mathbf{x})$ with:

$$K_{\text{issf}}(\mathbf{x}) \triangleq \left\{ \mathbf{u} \in \mathbb{R}^m \mid \begin{array}{l} \dot{h}_{\Pi}(\Pi(\mathbf{x}), \mathbf{u}) \geq \\ -\alpha(h_{\Pi}(\Pi(\mathbf{x}))) - \iota(\|\delta\|) \end{array} \right\},$$

Proof. As h_{Π} is an ISSf-CBF for (15) on \mathcal{C}_{Π} and the state-feedback controller satisfies $\mathbf{k}(\mathbf{x}) \in K_{\text{issf}}(\mathbf{x})$, Theorem 2 implies that the controller \mathbf{k} renders the set \mathcal{C}_{Π} input-to-state safe for all δ satisfying $\|\delta\|_{\infty} \leq \bar{\delta}$. In particular, there exists $\gamma \in \mathcal{K}_{\infty}$ such that the set:

$$\mathcal{C}_{\Pi, \delta} \triangleq \{ \mathbf{y} \in \mathbb{R}^k \mid h_{\Pi}(\mathbf{y}) + \gamma(\|\delta\|_{\infty}) \geq 0 \}, \quad (23)$$

is safe. Let $\mathbf{x}_0 \in \mathbb{R}^n$ be such that $\mathbf{y}_0 = \Pi(\mathbf{x}_0) \in \mathcal{C}_{\Pi, \delta}$. With $\mathbf{x}(0) = \mathbf{x}_0$ (implying $\mathbf{y}(0) = \mathbf{y}_0$), safety of $\mathcal{C}_{\Pi, \delta}$ implies:

$$h_{\Pi}(\Pi(\mathbf{x}(t))) + \gamma(\|\delta\|_{\infty}) \geq 0, \quad (24)$$

for $t \in I(\mathbf{x}_0)$. As h_{Π} is a compatible projection for h with respect to Π , we have:

$$\bar{\sigma}(h(\mathbf{x}(t))) + \gamma(\|\delta\|_{\infty}) \geq 0, \quad (25)$$

Multiplying both sides by $\frac{1}{2}$ and using that $\bar{\sigma} \in \mathcal{K}_{\infty, e}$, it follows that:

$$\bar{\sigma}^{-1} \left(\frac{1}{2} \bar{\sigma}(h(\mathbf{x}(t))) + \frac{1}{2} \gamma(\|\delta\|_{\infty}) \right) \geq 0, \quad (26)$$

The triangle inequality for class \mathcal{K} functions [12] implies:

$$h(\mathbf{x}(t)) + \underbrace{\bar{\sigma}^{-1}(\gamma(\|\delta\|_{\infty}))}_{\gamma'(\|\delta\|_{\infty})} \geq 0, \quad (27)$$

for all $t \in I(\mathbf{x}_0)$, implying the set \mathcal{C}_{δ} defined as in (16)-(18) using γ' is forward invariant, and hence safe. Thus the closed-loop system (7) is PSSf on \mathcal{C} with respect to Π and corresponding projected disturbances δ . \square

Corollary 1. Let $\mathcal{C} \subset \mathbb{R}^n$ be the 0-superlevel set of a continuously differentiable function $h : \mathbb{R}^n \rightarrow \mathbb{R}$ with 0 a regular value. Viewing h as a projection such that $y = h(\mathbf{x})$, let the projected dynamics be given by:

$$\dot{y} = f_y(\mathbf{x}) + \mathbf{g}_y(\mathbf{x})\mathbf{u} + \delta \quad (28)$$

with projected disturbances $\delta \in \mathbb{R}$. If there exists a Lipschitz continuous feedback controller $\mathbf{k} : \mathbb{R}^n \rightarrow \mathbb{R}^m$ such that:

$$f_y(\mathbf{x}) + \mathbf{g}_y(\mathbf{x})\mathbf{k}(\mathbf{x}) \geq -\alpha(y), \quad (29)$$

and there exists $\bar{\delta} > 0$ satisfying $|\delta| < \bar{\delta}$, then the disturbed system (7) can be rendered PSSf on \mathcal{C} with respect to the projection h and projected disturbances δ .

Proof. We first note that the identity map $I : \mathbb{R} \rightarrow \mathbb{R}$ is a compatible projection for h :

$$h(\mathbf{x}) \leq I(h(\mathbf{x})) \leq h(\mathbf{x}) \quad (30)$$

with $\underline{\sigma}(r) = \bar{\sigma}(r) = r$. Furthermore, the inequality in (29) implies the identity map can be viewed as an ISSf-CBF for the projected dynamics (28):

$$\sup_{\mathbf{u} \in \mathbb{R}^m} \dot{I}(\mathbf{x}, \mathbf{u}, \delta) \geq \dot{I}(\mathbf{x}, \mathbf{k}(\mathbf{x}), \delta) \geq -\alpha(I(y)) - |\delta|, \quad (31)$$

for all $\mathbf{x} \in \mathbb{R}^n$ and $\delta \in \mathbb{R}$ satisfying $|\delta| \leq \bar{\delta}$. Therefore the system (7) can be rendered PSSf on \mathcal{C} with respect to the projection h and projected disturbances δ by Theorem 3. \square

IV. INTEGRATION WITH LEARNING

In this section we consider a structured form of uncertainty in affine control systems. We discuss the impact of this uncertainty in a CBF time derivative, and on the PSSf behavior of the system. We demonstrate how learning can be used to mitigate the resulting impact on safety.

In practice, the system dynamics (1) are not known during control design due to parametric error and unmodeled dynamics. Instead, a nominal model of the system is utilized:

$$\dot{\hat{\mathbf{x}}} = \hat{\mathbf{f}}(\mathbf{x}) + \hat{\mathbf{g}}(\mathbf{x})\mathbf{u}, \quad (32)$$

where $\hat{\mathbf{f}} : \mathbb{R}^n \rightarrow \mathbb{R}^n$ and $\hat{\mathbf{g}} : \mathbb{R}^n \rightarrow \mathbb{R}^{n \times m}$ are assumed to be Lipschitz continuous on \mathbb{R}^n . By adding and subtracting (32) to (1), the dynamics of the system can be expressed as:

$$\dot{\mathbf{x}} = \hat{\mathbf{f}}(\mathbf{x}) + \hat{\mathbf{g}}(\mathbf{x})\mathbf{u} + \underbrace{\mathbf{f}(\mathbf{x}) - \hat{\mathbf{f}}(\mathbf{x})}_{\mathbf{b}(\mathbf{x})} + \underbrace{(\mathbf{g}(\mathbf{x}) - \hat{\mathbf{g}}(\mathbf{x}))\mathbf{u}}_{\mathbf{A}(\mathbf{x})}, \quad (33)$$

where the unknown disturbance $\mathbf{d} = \mathbf{b}(\mathbf{x}) + \mathbf{A}(\mathbf{x})\mathbf{u}$ is assumed to be time invariant, but explicitly depends on the state and input to the system.

If the function $h : \mathbb{R}^n \rightarrow \mathbb{R}$ is a CBF for the nominal model (32) on \mathcal{C} , the uncertainty in the dynamics directly manifests in the time derivative of h :

$$\begin{aligned} \dot{h}(\mathbf{x}, \mathbf{u}) &= \overbrace{\frac{\partial h}{\partial \mathbf{x}}(\mathbf{x}) (\hat{\mathbf{f}}(\mathbf{x}) + \hat{\mathbf{g}}(\mathbf{x})\mathbf{u})}^{\hat{h}(\mathbf{x}, \mathbf{u})} \\ &+ \underbrace{\frac{\partial h}{\partial \mathbf{x}}(\mathbf{x}) \mathbf{b}(\mathbf{x})}_{\mathbf{b}(\mathbf{x})} + \underbrace{\frac{\partial h}{\partial \mathbf{x}}(\mathbf{x}) \mathbf{A}(\mathbf{x}) \mathbf{u}}_{\mathbf{a}(\mathbf{x})^{\top}}. \end{aligned} \quad (34)$$

Given that h is a CBF for (32) on \mathcal{C} , let $\mathbf{k} : \mathbb{R}^n \rightarrow \mathbb{R}^m$ be a Lipschitz continuous state-feedback controller such that:

$$\sup_{\mathbf{u} \in \mathbb{R}^m} \hat{h}(\mathbf{x}, \mathbf{u}) \geq \hat{h}(\mathbf{x}, \mathbf{k}(\mathbf{x})) \geq -\alpha(h(\mathbf{x})). \quad (35)$$

Letting the projected disturbance be defined as:

$$\delta = \mathbf{b}(\mathbf{x}) + \mathbf{a}(\mathbf{x})^{\top} \mathbf{k}(\mathbf{x}), \quad (36)$$

Corollary 1 implies that if there exists a $\bar{\delta} > 0$ such that $|\mathbf{b}(\mathbf{x}) + \mathbf{a}(\mathbf{x})^{\top} \mathbf{k}(\mathbf{x})| \leq \bar{\delta}$ for all $\mathbf{x} \in \mathbb{R}^n$, the uncertain system (1) can be rendered PSSf on \mathcal{C} with respect to the projection h and projected disturbances δ .

As in [20], we may wish to reduce the error between \dot{h} and \hat{h} by utilizing data-driven models to estimate the functions

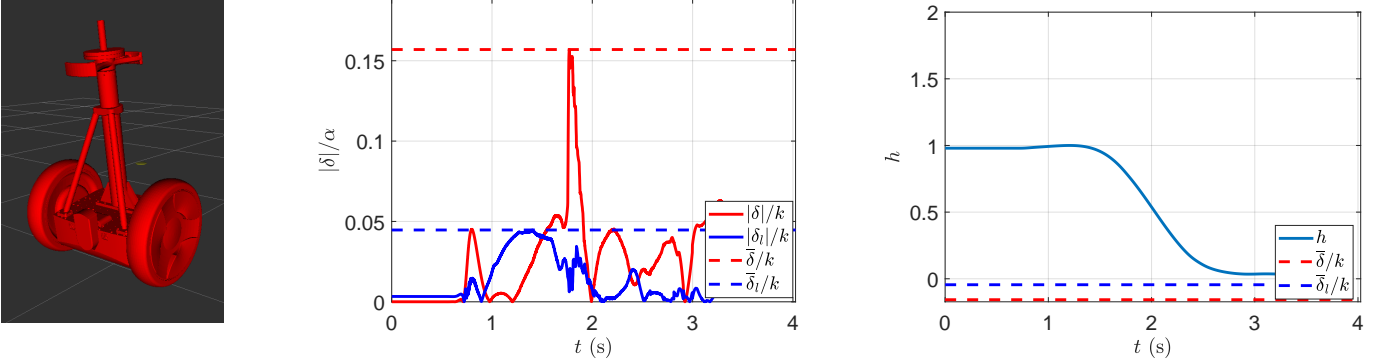


Fig. 1. Simulation results with Segway platform demonstrating improvement in PSSf behavior. **(Left)** Robotic Segway platform model used in simulation. **(Center)** Absolute value of the projected disturbance δ along the trajectory without learning models ((36),red) and with learning models ((38), blue), with learning reducing the worse case projected disturbance ($\bar{\delta}/\alpha$). **(Right)** The value of the barrier satisfies the corresponding worst case lower bound with and without learning being used to compute δ . The worst case lower bound is raised with learning (the blue dashed line lies above the red dashed line).

b and a . In particular, given Lipschitz continuous estimators $\hat{b} : \mathbb{R}^n \rightarrow \mathbb{R}$ and $\hat{a} : \mathbb{R}^n \rightarrow \mathbb{R}^m$, (34) can be reformulated as:

$$\begin{aligned} \dot{h}(\mathbf{x}, \mathbf{u}) = & \underbrace{\frac{\partial h}{\partial \mathbf{x}}(\mathbf{x}) \left(\hat{\mathbf{f}}(\mathbf{x}) + \hat{\mathbf{g}}(\mathbf{x})\mathbf{u} \right)}_{\hat{h}(\mathbf{x}, \mathbf{u})} + \hat{b}(\mathbf{x}) + \hat{a}(\mathbf{x})^\top \mathbf{u} \\ & + \underbrace{\frac{\partial h}{\partial \mathbf{x}}(\mathbf{x}) \mathbf{b}(\mathbf{x}) - \hat{b}(\mathbf{x})}_{\tilde{b}(\mathbf{x})} + \underbrace{\left(\frac{\partial h}{\partial \mathbf{x}}(\mathbf{x}) \mathbf{A}(\mathbf{x}) - \hat{a}(\mathbf{x})^\top \right)}_{\tilde{a}(\mathbf{x})^\top} \mathbf{u}. \end{aligned} \quad (37)$$

Under the assumption that the introduction of the estimators does not violate the CBF condition, such that there exists a state-feedback controller \mathbf{k} satisfying (35) with \hat{h} defined as in (37), we may define the projected disturbance as:

$$\delta = \tilde{b}(\mathbf{x}) + \tilde{a}(\mathbf{x})^\top \mathbf{k}(\mathbf{x}) \quad (38)$$

As before, if there exists $\bar{\delta} > 0$ such that $|\tilde{b}(\mathbf{x}) + \tilde{a}(\mathbf{x})^\top \mathbf{k}(\mathbf{x})| \leq \bar{\delta}$ for all $\mathbf{x} \in \mathbb{R}^n$, Corollary 1 can be used to certify (1) as PSSf on \mathcal{C} with respect to the projection h and projected disturbances δ . The preceding statements are formalized in the following theorem:

Theorem 4. *Let $\mathcal{C} \subset \mathbb{R}^n$ be the 0-superlevel set of a continuously differentiable function $h : \mathbb{R}^n \rightarrow \mathbb{R}$ with 0 a regular value, and let $\hat{h} : \mathbb{R}^n \rightarrow \mathbb{R}$ be defined as in (34) or (37). If there exists a Lipschitz continuous state-feedback controller $\mathbf{k} : \mathbb{R}^n \rightarrow \mathbb{R}^m$ satisfying (35), and $\bar{\delta} > 0$ such that the corresponding projected disturbance defined as in (36) or (38) satisfies $|\delta| \leq \bar{\delta}$, then (1) is PSSf on \mathcal{C} with respect to the projection h and projected disturbances δ .*

In the presence of estimators, this theorem defines a quantitative relationship between the prediction error of the estimators, $|\dot{h}(\mathbf{x}, \mathbf{k}(\mathbf{x})) - \hat{h}(\mathbf{x}, \mathbf{k}(\mathbf{x}))| = |\delta|$, and the degradation of the safety of the closed-loop system. As the prediction error is reduced (via additional training data or more complex learning models), the set kept safe more closely resembles \mathcal{C} .

V. EXPERIMENTAL VALIDATION

To demonstrate the ability of learning to improve safety guarantees via Projection-to-State Safety, we deployed the episodic learning framework with CBFs established in [20] on a robotic Segway platform, seen in Figure 1 and 2, in simulation and experimentally. The planar, 4 dimensional, Segway was considered, with states given by horizontal position, horizontal velocity, pitch angle, and pitch angle rate. The input the system is specified as a torque about the wheel at the base of the Segway. In both cases a sequence of episodes were ran to train estimators \hat{b} and \hat{a} .

In each episode the Segway was set to track a desired trajectory in the pitch angle space without violating a barrier function on a portion of its state, using the safety-critical control formulation in [11]. After the sequence of episodes, the Segway was ran once more with a learning-informed controller, and the projected disturbance δ as defined in (36) and (38) was computed. The worst case disturbance $\bar{\delta}$ was found, and a lower bound on h for that trajectory was determined using the fact $h \leq \alpha^{-1}(\bar{\delta}) \implies \dot{h} \geq 0$. In both simulation and experimental results, $\alpha(r) = kr$ with $k > 0$.

In simulation, the Segway was given a bound on its position in space, forcing it to stay within one meter of its starting location. The CBF was generated through the backup controller method [10]. The value of the CBF is computed at each time-step by integrating the system forward in time under a backup control law. Sensitivity analysis along the trajectory is used to compute the gradient of the CBF. This simulation result highlights the ability of learning to reduce worst case disturbances for complex CBFs that cannot be expressed in closed-form. The simulation was done in a ROS-based C++ environment [18]. The simulation environment accurately simulates the physical system by adding input delay, sensor noise, and state estimation. Experimentally, a simple CBF was specified to limit the pitch angle and pitch angle rate of the Segway to an ellipse about the Segway's equilibrium state. The desired pitch angle trajectory would lead to the Segway tipping quickly, thereby violating the safety set in the absence of the CBF and safety-critical control formulation.

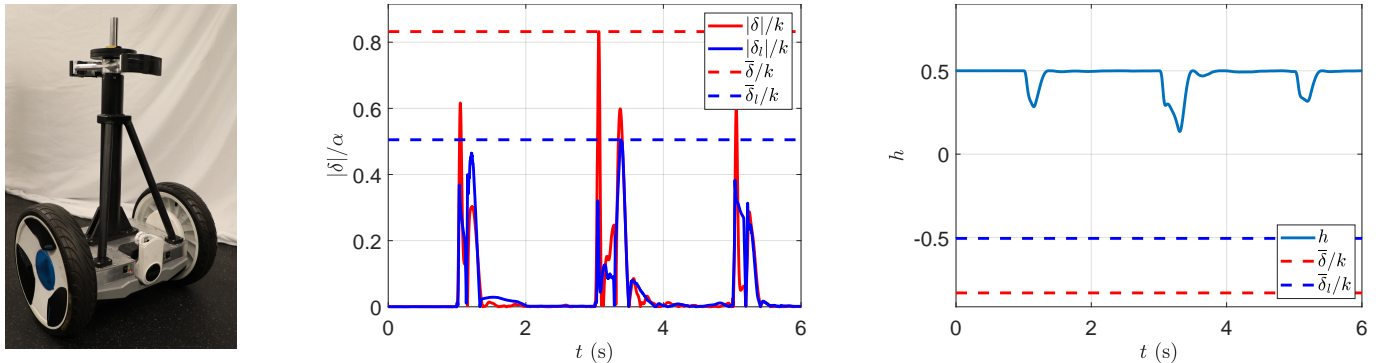


Fig. 2. Experimental results with Segway platform demonstrating improvement in PSSf behavior. **(Left)** Physical robotic Segway platform used in experimentation. **(Center)** Absolute value of the projected disturbance δ along the trajectory without learning models ((36), red) and with learning models ((38), blue), with learning reducing the worst case projected disturbance ($\bar{\delta}/\alpha$). **(Right)** The value of the barrier satisfies the corresponding worst case lower bound with and without learning being used to compute δ . The worst case lower bound is raised with learning (the blue dashed line lies above the red dashed line).

In both cases, we see that introducing learning estimators into the computation of the projected disturbance decreases the worst case disturbance ($\bar{\delta} > \bar{\delta}_l$). This leads to a greater lower bound on h , and thus a stronger guarantee on the PSSf behavior of the system. We note that the conservative nature of the lower bounds on h arise from the fact that the worst case disturbance $\bar{\delta}$ along the trajectory is used. If the worst case disturbance can be reduced (by data-aware control synthesis), stronger guarantees on safety can be made.

VI. CONCLUSIONS

We presented a novel method for assessing the impact of disturbances on safety in a project environment via Projection-to-State Safety, and considered how it can be utilized in conjunction with learning to mitigate the impact of model uncertainty on safety. We demonstrate the ability of learning to improve the guarantees endowed by PSSf in simulation and experimentally on a Segway platform. Future work includes developing data-driven methods for quantifying the worst case projected disturbance, and synthesizing data-aware controllers that reduce the projected disturbance.

REFERENCES

- [1] A. Ames, J. Grizzle, and P. Tabuada. Control barrier function based quadratic programs with application to adaptive cruise control. In *Conference on Decision & Control (CDC)*, pages 6271–6278. IEEE, 2014.
- [2] A. D. Ames, S. Coogan, M. Egerstedt, G. Notomista, K. Sreenath, and P. Tabuada. Control barrier functions: Theory and applications. In *European Control Conference (ECC)*, pages 3420–3431. IEEE, 2019.
- [3] A. D. Ames, P. Tabuada, and S. Sastry. On the stability of zeno equilibria. In *International Workshop on Hybrid Systems: Computation and Control (HSCC)*, pages 34–48. Springer, 2006.
- [4] A. D. Ames, X. Xu, J. W. Grizzle, and P. Tabuada. Control barrier function based quadratic programs for safety critical systems. *IEEE Transactions on Automatic Control*, 62(8):3861–3876, 2017.
- [5] F. Berkenkamp, A. P. Schoellig, and A. Krause. Safe controller optimization for quadrotors with gaussian processes. In *International Conference on Robotics and Automation (ICRA)*, pages 491–496. IEEE, 2016.
- [6] F. Blanchini and S. Miani. *Set-theoretic methods in control*. Springer, 2008.
- [7] X. Chen, M. Monfort, A. Liu, and B. D. Ziebart. Robust covariate shift regression. In *Artificial Intelligence and Statistics (AISTATS)*, pages 1270–1279, 2016.
- [8] R. Cheng, G. Orosz, R. M. Murray, and J. W. Burdick. End-to-end safe reinforcement learning through barrier functions for safety-critical continuous control tasks. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 33, pages 3387–3395, 2019.
- [9] J. F. Fisac, A. K. Akametalu, M. N. Zeilinger, S. Kaynama, J. Gillula, and C. J. Tomlin. A general safety framework for learning-based control in uncertain robotic systems. *IEEE Transactions on Automatic Control*, 2018.
- [10] T. Gurriet, M. Mote, A. D. Ames, and E. Feron. An online approach to active set invariance. In *Conference on Decision and Control (CDC)*, pages 3592–3599. IEEE, 2018.
- [11] T. Gurriet, A. Singletary, J. Reher, L. Ciarletta, E. Feron, and A. Ames. Towards a framework for realizable safety critical control through active set invariance. In *International Conference on Cyber-Physical Systems*, pages 98–106. IEEE Press, 2018.
- [12] C. M. Kellett. A compendium of comparison function results. *Mathematics of Control, Signals, and Systems*, 26(3):339–374, 2014.
- [13] S. Kolathaya and A. D. Ames. Input-to-state safety with control barrier functions. *IEEE Control Systems Letters*, 3(1):108–113, 2018.
- [14] A. Liu, G. Shi, S.-J. Chung, A. Anandkumar, and Y. Yue. Robust regression for safe exploration in control. In *Conference on Learning for Decision and Control (LADC)*, 2020.
- [15] Q. Nguyen and K. Sreenath. Exponential control barrier functions for enforcing high relative-degree safety-critical constraints. In *American Control Conference (ACC)*, pages 322–328. IEEE, 2016.
- [16] M. Ohnishi, L. Wang, G. Notomista, and M. Egerstedt. Barrier-certified adaptive reinforcement learning with applications to brushbot navigation. *IEEE Transactions on Robotics*, 35(5):1186–1205, 2019.
- [17] L. Perko. *Differential equations and dynamical systems*, volume 7. Springer Science & Business Media, 2013.
- [18] A. Singletary, P. Nilsson, T. Gurriet, and A. D. Ames. Online active safety for robotic manipulators. In *International Conference on Intelligent Robots and Systems (IROS)*, pages 173–178. IEEE, 2019.
- [19] A. J. Taylor, V. D. Dorobantu, M. Krishnamoorthy, H. M. Le, Y. Yue, and A. D. Ames. A control lyapunov perspective on episodic learning via projection to state stability. In *Conference on Decision & Control (CDC)*, pages 1448–1455. IEEE, 2019.
- [20] A. J. Taylor, A. Singletary, Y. Yue, and A. D. Ames. Learning for safety-critical control with control barrier functions. In *Conference on Learning for Decision and Control (LADC)*, 2020.
- [21] L. Wang, A. D. Ames, and M. Egerstedt. Safety barrier certificates for collisions-free multirobot systems. *IEEE Transactions on Robotics*, 33(3):661–674, 2017.
- [22] L. Wang, E. A. Theodorou, and M. Egerstedt. Safe learning of quadrotor dynamics using barrier certificates. In *International Conference on Robotics and Automation (ICRA)*, pages 2460–2465. IEEE, 2018.
- [23] X. Xu, P. Tabuada, J. W. Grizzle, and A. D. Ames. Robustness of control barrier functions for safety critical control. *IFAC-PapersOnLine*, 48(27):54–61, 2015.