

Formal Verification of Safety Critical Autonomous Systems via Bayesian Optimization

Prithvi Akella, Ugo Rosolia, Andrew Singletary, and Aaron D. Ames¹

Abstract—As control systems become increasingly more complex, there exists a pressing need to find systematic ways of verifying them. To address this concern, there has been significant work in developing test generation schemes for black-box control architectures. These schemes test a black-box control architecture’s ability to satisfy its control objectives, when these objectives are expressed as operational specifications through temporal logic formulae. Our work extends these prior, model based results by lower bounding the probability by which the black-box system will satisfy its operational specification, when subject to a pre-specified set of environmental phenomena. We do so by systematically generating tests to minimize a Lipschitz continuous robustness measure for the operational specification. We demonstrate our method with experimental results, wherein we show that our framework can reasonably lower bound the probability of specification satisfaction.

I. INTRODUCTION

An integral aspect of control system design is to ensure that the driven system satisfies some extraneous criteria, *e.g.* safety, robustness to noise/model mismatch, *etc.* To quantify this assurance, these criteria are oftentimes expressed as temporal logic formulae, and the development of controllers which are guaranteed to satisfy these formulae - termed correct-by-construction controllers - have seen significant interest in the recent past [2], [3]. However, as control systems become increasingly more complex and are subject to more diverse, varied scenarios, developing a correct-by-construction controller becomes progressively more difficult, if not, perhaps, impossible [4]. As a result, the problem remains, as to how to verify that a complex control architecture satisfies the criteria required of it.

In the Test and Evaluation (T&E) community, significant work has been done to address this verification dilemma. Specifically, there has been some work to extend traditional safety verification techniques, by iteratively solving for candidate Lyapunov/Barrier certificates based on simulation data/a dynamic model [5]–[7]. However, the ideal would be to verify the controller’s ability to satisfy its criteria despite the presence of confounding environmental phenomena, as these oftentimes lead to a system’s inability to satisfy specifications *e.g.* unpredictable human behavior for autonomous cars and adversarial agents in a reconnaissance context [4].

Numerous, model-based approaches to testing for verification have been studied by the T&E community [8]–[10]; however, they tend to be sample inefficient. As a result, there

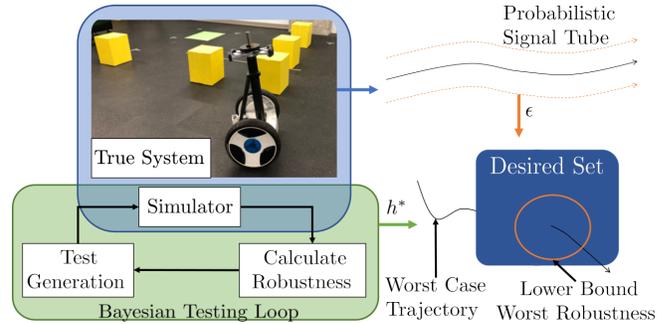


Fig. 1. This paper presents a method to lower bound the probability with which a black-box system (the Segway shown top left) satisfies its operational specification (please see Figures 2 and 3 for the full experimental setup). This method first determines the worst-case robustness measure, h^* , for a system simulator, and uses this h^* and a norm-bound, ϵ , from the signal tube, to lower bound the worst-case robustness for the black-box system, to some minimum probability.

has also been significant work in increasing sample efficiency of these methods [11], [12]. One such sample efficient method is Bayesian Optimization, which has seen success in terms of fine tuning existing control strategies [13], [14]. From a testing perspective, identifying environmental phenomena to frustrate satisfaction of the system’s operational specifications can also be recast to a minimization problem solvable by Bayesian Optimization [15]. Indeed, there has also been work to identify multiple such phenomena, should they exist [16].

Our Contribution We extend prior results in the Bayesian testing community, by lower bounding the probability with which the actual system satisfies its operational specifications. We do so for Reach-Avoid type STL specifications, as they permit a Lipschitz continuous robustness measure. As a result, the testing for verification problem can be reformulated as an optimization problem that is solvable by Bayesian Optimization - minimizing this robustness measure subject to an adversarial environment. Finally, we use the solution to this minimization problem to lower bound the probability by which the true system satisfies the same operational specification.

Outline in Section II-A, we provide a brief mathematical overview of relevant topics. Section II-B formally sets up the problem under study. Section III-A details the main theorem and its proof, and Section III-B states and proves three propositions generalizing the theorem. Finally, Section IV shows how our method provides a reasonable lower bound on the probability that a Segway satisfies its safety specification, by verifying its associated controller.

* This work was supported by the Air Force Office of Scientific Research.

¹ The authors are with the California Institute of Technology, 1200 East California Boulevard, Pasadena, CA 91125, USA. pakella@caltech.edu, urosolia@caltech.edu, asinglet@caltech.edu, ames@caltech.edu

II. PROBLEM FORMULATION

This section details the necessary mathematical information for the sequel and frames the problem under study.

A. Mathematical Preliminaries

This subsection will be split into three parts encompassing some general notation, a brief description on Signal Temporal Logic, and a brief description of Bayesian Optimization.

Notation: \mathbb{R}^d denotes the d -dimensional Euclidean Space. A function, $f : \mathbb{X} \rightarrow \mathbb{Y}$, is $(L, \|\cdot\|_{\mathbb{X}}, \|\cdot\|_{\mathbb{Y}})$ -Lipschitz continuous, if there exists a positive constant, L , such that: $\|f(x) - f(y)\|_{\mathbb{Y}} \leq L\|x - y\|_{\mathbb{X}} \forall x, y \in \mathbb{X}$, where $\|\cdot\|_{\mathbb{X}}$ and $\|\cdot\|_{\mathbb{Y}}$ are norms for the normed vector spaces, \mathbb{X} and \mathbb{Y} , respectively. A signal is a function, $s : \mathbb{R}_+ \rightarrow \mathbb{X}$, that maps time to a vector in a vector space, \mathbb{X} , i.e. $s(t) = x \in \mathbb{X}$. $\mathcal{S}_{t_f}^{\mathbb{X}}$ is the vector space of all signals over the bounded time-frame, $[0, t_f]$, i.e. $\mathcal{S}_{t_f}^{\mathbb{X}} = \{s : \mathbb{R}_+ \rightarrow \mathbb{X} \mid \exists x \in \mathbb{X} \text{ s.t. } s(t) = x \forall t \in [0, t_f]\}$. If the vector space, \mathbb{X} , admits a norm, $\|\cdot\|$, $\|\cdot\|^{t_f} = \max_{t \in [0, t_f]} \|s(t)\|$, is its induced norm over $\mathcal{S}_{t_f}^{\mathbb{X}}$.

Signal Temporal Logic: Signal Temporal Logic is a language by which rich, time-varying system behavior can be succinctly expressed. This language is based on atomic propositions, $\Phi \in \mathcal{A}$, which are boolean valued variables dependent on system behavior:

$$\Phi(x) = \text{True} \iff x \in \llbracket \Phi \rrbracket = \{x \in \mathbb{R}^n \mid b(x) \sim \mu\}. \quad (1)$$

Here, \mathcal{A} is the set of all atomic propositions, $\Phi(x)$ denotes the truth evaluation of Φ at the state, $x, b : \mathbb{R}^n \rightarrow \mathbb{R}, \mu \in \mathbb{R}$, and $\sim = \{\geq, \leq, <, >\}$ [17]–[19]. Additionally, \mathcal{A} is closed under logical combinations of its components, i.e.:

$$\Phi \in \mathcal{A} \implies \neg\Phi \in \mathcal{A}, \quad (2)$$

$$\Phi_1, \Phi_2 \in \mathcal{A} \implies \Phi_1 \wedge \Phi_2 \in \mathcal{A} \text{ and } \Phi_1 \vee \Phi_2 \in \mathcal{A}. \quad (3)$$

Here, \neg denotes negation, \wedge denotes conjunction (and), and \vee denotes disjunction (or). System specifications, ψ , can be defined as follows:

$$\psi \triangleq \text{True} \mid \Phi \mid \neg\psi \mid \psi_1 \vee \psi_2 \mid \psi_1 \wedge \psi_2 \mid \psi_1 \text{ U } \psi_2, \quad (4)$$

where, ψ_1, ψ_2 are specifications themselves [20]. Finally, \mathbb{S} is the set of all STL specifications, i.e. $\psi \in \mathbb{S}$ [20].

When a signal, $s : \mathbb{R}_+ \rightarrow \mathbb{R}^n$, satisfies a specification, ψ , by some time, t i.e. ensures $\psi = \text{True}$, we write $s(t) \models \psi$. Here, \models is termed the satisfaction relation, and it is defined as follows:

$$s(t) \models \Phi \iff \Phi(s(t)) = \text{True}, \quad (5)$$

$$s(t) \models \neg\psi \iff \psi(s(t_f)) = \text{False}, \quad (6)$$

$$s(t) \models \psi_1 \vee \psi_2 \iff s(t) \models \psi_1 \vee s(t) \models \psi_2, \quad (7)$$

$$s(t) \models \psi_1 \wedge \psi_2 \iff s(t) \models \psi_1 \wedge s(t) \models \psi_2, \quad (8)$$

$$s(t) \models \psi_1 \text{ U } \psi_2 \iff \exists t^* \leq t \text{ s.t.} \quad (9)$$

$$(s(t') \models \psi_1 \forall t' < t^*) \wedge (s(t^*) \models \psi_2). \quad (10)$$

Finally, for any signal temporal logic specification, ψ , evaluated over some bounded time-frame, $[0, t]$, there exists a

robustness measure, ρ , with which to measure proximity to satisfaction of the specification [17]:

$$\rho : \mathcal{S}_t^{\mathbb{R}^n} \rightarrow \mathbb{R} \text{ s.t. } \rho(s) \geq 0 \iff s(t) \models \psi. \quad (11)$$

Bayesian Optimization: Mathematically, Bayesian Optimization is a solution procedure intended to solve optimization problems of the following form:

$$x^* = \underset{x \in A \subseteq \mathbb{R}^d}{\text{argmin}} c(x), \quad (12)$$

where, typically, d is small; $c(x)$ is expensive to evaluate and lacks nice, analytic structure e.g. convexity; and A is some set for which membership evaluation is simple e.g. a hyperrectangle [21].

Abstractly, the procedure follows two, main steps, (for a more comprehensive mathematical treatment, please reference [22]). The first step fits a Gaussian Process to c based on a data-set of sampled values, $\mathcal{D}_N = \{(x_k, y_k = c(x_k))\}_{k=1}^N$ and a set of parameters, θ . The second step finds the next, sample point via optimizing an acquisition function - in our case, the Expected Improvement function. Usually, the parameters θ are updated after each cycle; however, the convergence results of [22] require a static parameterization. Therefore, we opt for a static parameterization.

B. Problem Setup

We consider an uncertain control system as follows:

$$\dot{x} = f(x, u, d, w), \quad (13)$$

where the state $x \in \mathcal{X} \subseteq \mathbb{R}^n$, the control input, $u \in \mathcal{U} \subseteq \mathbb{R}^m$, the environmental configuration, $d \in \mathcal{D} \subseteq \mathbb{R}^p$, and the disturbance, $w \sim \pi_{\text{env}}$, for the unknown distribution, π_{env} . For this system, we also have a controller,

$$U : \mathbb{R}^n \times \mathbb{R}^p \rightarrow \mathcal{U} \subseteq \mathbb{R}^m, \text{ s.t. } u(t) = U(x(t), d). \quad (14)$$

Our goal is to certify that the above controller, U , which is designed based on a nominal model,

$$\dot{x} = \hat{f}(x, u, d), \quad (15)$$

satisfies a specification, ψ , to some minimum probability, for all environment configurations, $d \in \mathcal{D}$, when in closed-loop with the unknown, uncertain system (13).

For the true and nominal systems from (13) and (15), we define a simulation trajectory, $\hat{\phi}(\cdot)$, and the actual system trajectory it approximates, $\phi(\cdot)$, as follows:

$$\hat{\phi}(x_0, u(t), d, t) = x_0 + \int_0^t \hat{f}(x(s), u(s), d) ds,$$

$$\phi(x_0, u(t), d, w(t), t) = x_0 + \int_0^t f(x(s), u(s), d, w(s)) ds.$$

As stated, both $\phi(x_0, u(t), d, w(t))$ and $\hat{\phi}(x_0, u(t), d)$ are signals, under the assumption that neither system ever engenders state values that tend to infinity.

As the simulation model (15) is deterministic, it may be different from the true system (13). Therefore, we introduce the following function to quantify the difference between solutions to the nominal, closed loop system,

$\hat{\phi}(x_0, U(x(t), d), d, t)$, and the true, closed loop system, $\phi(x_0, U(x(t), d), d, w(t), t)$:

$$\begin{aligned} \Delta_\phi(x_0, U, d, w, t) \\ = \phi(x_0, U(x(t), d), d, w(t), t) - \hat{\phi}(x_0, U(x(t), d), d, t), \end{aligned} \quad (17)$$

where $w(t) \in \mathcal{W}_t$ is a random signal where, $w(t) \sim \pi_{\text{env}} \forall t \in [0, t_f]$, and \mathcal{W}_{t_f} is the set of all random signals over the bounded time-frame, $[0, t_f]$. Equation (17) permits us to define the accuracy of our nominal closed loop system (14)-(15), over bounded time-frames, $[0, t_f]$:

Definition 1. The nominal closed loop system (14)-(15) is $(\epsilon, t_f, \lambda, \|\cdot\|)$ -accurate, if and only if

$$\mathbb{P}_w \left[\max_{0 \leq t \leq t_f} \|\Delta_\phi(x_0, U, d, w(t), t)\| \leq \epsilon \right] \geq 1 - \lambda, \quad (18)$$

$\forall d \in \mathcal{D}$ where $\Delta_\phi(x_0, U, d, w)$ is defined in (17).

We require a notion of accuracy for our nominal model, as we will first certify that the nominal closed loop system (14)-(15) satisfies the specification, ψ , for all environmental configurations, $d \in \mathcal{D}$. If the nominal closed loop system indeed satisfies ψ , then, provided that ψ meets the criteria for the following assumption, we will extend that certificate to the actual system, with some minimum probability:

Assumption 1. The specification, ψ , is of the form, $\psi = \text{True} \cup \Phi$ or $\psi = \neg(\text{True} \cup \neg\Phi)$, where $\Phi \in \mathcal{A}$.

In effect then, our problem statement is as follows.

Problem Statement. For a specification, ψ , satisfying the criteria for Assumption 1, identify the minimum probability by which the true closed loop system (13)-(14) satisfies ψ by some final time, t_f , i.e. determine p where,

$$\mathbb{P}_w [\phi(x_0, U(x(t), d), d, w(t), t_f) \models \psi] \geq p, \quad \forall d \in \mathcal{D}, \quad (19)$$

III. MAIN RESULT

This section is divided into two parts. First, we state and prove the main result. Then, we show under which conditions the assumption from our main result are satisfied.

A. Statement of Main Result

When a system's operational specification, ψ , satisfies the criteria for Assumption 1, and the controller (14) is designed based on the nominal model (15), we devise a method to lower bound the probability by which the true closed loop system (13)-(14) satisfies ψ , for all environmental configurations. We do so by first assuming that there always exists a Lipschitz robustness measure, $\rho : \mathcal{S}_{t_f}^{\mathbb{R}^n} \rightarrow \mathbb{R}$, for specifications, ψ , satisfying Assumption 1 (we will prove that such a Lipschitz measure always exists for these specifications, in Section III-B). Then, we define a nominal, worst-case robustness value, h^* , for this robustness measure, ρ , and the nominal closed loop system (14)-(15):

$$h^* = \min_{d \in \mathcal{D}} \rho(\hat{\phi}(x_0, u(t), d)). \quad (20)$$

Finally, our result indicates that if h^* is both positive and sufficiently far from zero, then the true closed loop system (13)-(14) is guaranteed to satisfy ψ to the same minimum probability defined by our simulator 1.

Theorem 1. Let ψ be a specification satisfying Assumption 1 with an $(L, \|\cdot\|^{t_f}, |\cdot|)$ -Lipschitz continuous robustness measure, $\rho : \mathcal{S}_{t_f}^{\mathbb{R}^n} \rightarrow \mathbb{R}$ as defined in (11). If the nominal system (14)-(15) is $(\epsilon, t_f, \lambda, \|\cdot\|)$ -accurate, and the nominal worst-case robustness, $h^* \geq L\epsilon$, then the true system (13)-(14) satisfies ψ with at least probability $1 - \lambda$, $\forall d \in \mathcal{D}$, i.e.

$$\mathbb{P}_w [\phi(x_0, U(x(t), d), d, w(t), t_f) \models \psi] \geq 1 - \lambda, \quad \forall d \in \mathcal{D}.$$

Theorem 1 states that if the nominal system robustly satisfies ψ , $h^* \geq L\epsilon$, and is $(\epsilon, t_f, \lambda, \|\cdot\|)$ -accurate, then with probability greater than $1 - \lambda$, the true system will satisfy the same specification, ψ , as well.

Proof: We start by showing that if $h^* \geq 0$, then,

$$\hat{\phi}(x_0, U(x(t), d), d, t_f) \models \psi \quad \forall d \in \mathcal{D}. \quad (21)$$

This stems directly, as h^* is defined as the minimum value of the robustness measure, ρ , over all environmental configurations, $d \in \mathcal{D}$. More formally, we have that

$$h^* \geq 0 \implies \rho(\hat{\phi}(x_0, U(x(t), d), d)) \geq 0, \quad \forall d \in \mathcal{D}, \quad (22)$$

$$\equiv \hat{\phi}(x_0, U(x(t), d), d, t_f) \models \psi \quad \forall d \in \mathcal{D}. \quad (23)$$

Similar to equation (20), were we to define the true worst-case robustness measure, H^* , for the true system (13)-(14) and some $w(t) \in \mathcal{W}_{t_f}$,

$$H^* = \min_{d \in \mathcal{D}} \rho(\phi(x_0, U(x(t), d), d, w(t))), \quad (24)$$

then if $H^* \geq 0$, the true closed loop system (13)-(14), is guaranteed to satisfy ψ , $\forall d \in \mathcal{D}$ and this, specific $w(t)$. The remainder of the proof uses h^* to lower bound H^* for all possible $w(t) \in \mathcal{W}_{t_f}$ to at least some minimum probability. To simplify notation in its presentation, we abbreviate $\hat{\phi}(x_0, U(x(t), d), d) = \bar{\phi}$ and $\phi(x_0, U(x(t), d), d, w(t)) = \tilde{\phi}$:

$$h^* - H^* = \min_{d \in \mathcal{D}} \rho(\bar{\phi}) - \min_{d \in \mathcal{D}} \rho(\tilde{\phi}), \quad (25)$$

$$= \max_{d \in \mathcal{D}} -\rho(\tilde{\phi}) - \max_{d \in \mathcal{D}} -\rho(\bar{\phi}), \quad (26)$$

$$\leq \max_{d \in \mathcal{D}} |\rho(\tilde{\phi}) - \rho(\bar{\phi})|, \quad (27)$$

$$\leq L \max_{d \in \mathcal{D}} \max_{0 \leq t \leq t_f} \|\tilde{\phi}(t) - \bar{\phi}(t)\|, \quad (28)$$

$$\leq L\epsilon \text{ with at least probability } 1 - \lambda. \quad (29)$$

Note that the second to last line is valid $\forall w(t) \in \mathcal{W}_{t_f}$. The transition to the last line arises as the nominal system is $(\epsilon, t_f, \lambda, \|\cdot\|)$ -accurate. Furthermore, this probability is over the random signals, $w(t) \in \mathcal{W}_{t_f}$, which is why the random signal, $w(t)$, does not appear in the last line. Hence,

$$h^* \geq L\epsilon \implies \mathbb{P}_w [H^* \geq 0] \geq 1 - \lambda \quad \forall d \in \mathcal{D}, \quad (30)$$

$$\equiv \mathbb{P}_w [\tilde{\phi}(t_f) \models \psi] \geq 1 - \lambda \quad \forall d \in \mathcal{D}, \quad (31)$$

where we abbreviate $\phi(x_0, U(x(t), d), d, w(t)) = \tilde{\phi}$. ■

B. Extension of Main Result

Theorem 1 assumes existence of an $(L, \|\cdot\|^{t_f}, |\cdot|)$ -Lipschitz continuous robustness measure, ρ , for specifications, ψ , satisfying Assumption 1. While this seems restrictive, the following three propositions show that such a Lipschitz robustness measure, ρ , always exists for specifications, ψ , satisfying Assumption 1.

The first proposition develops a Lipschitz continuous function, h , with which to measure satisfaction of any predicate, $\Phi \in \mathcal{A}$:

Proposition 1. *Let $\Phi \in \mathcal{A}$. For some norm, $\|\cdot\|$ on \mathbb{R}^n , there exists an $(L, \|\cdot\|, |\cdot|)$ -Lipschitz function, $h : \mathbb{R}^n \rightarrow \mathbb{R}$, such that $h(x) \geq 0 \iff \Phi(x) = \text{True}$.*

Proof: For $\Phi = \text{True}$, any constant, positive function h suffices e.g. $h(x) = 1$; likewise, for $\Phi = \text{False}$, any constant, negative function h suffices. It remains to show the same result for each $\Phi \in \mathcal{A} \setminus \{\text{True}, \text{False}\}$. Here, we note that as $\Phi \neq \text{True}, \text{False}$, its associated truth region, $\llbracket \Phi \rrbracket$, has a non-trivial boundary, i.e. $\partial \llbracket \Phi \rrbracket \neq \emptyset$. As a result, we can define a signed, distance function, h ,

$$l(x) = \min_{y \in \partial \llbracket \Phi \rrbracket} \|x - y\|, \quad h(x) = \begin{cases} l(x) & x \in \llbracket \Phi \rrbracket, \\ -l(x) & x \in \llbracket \neg \Phi \rrbracket. \end{cases}$$

By definition, h is positive if and only if $\Phi = \text{True}$. Also, h carries the same Lipschitz constant, $L = 1$, as the set-distance function, l , thus completing the proof. ■

Similar to Proposition 1, the second proposition develops a Lipschitz robustness measure, ρ , for any specification, ψ , that satisfies Assumption 1:

Proposition 2. *Let ψ be a specification that satisfies Assumption 1, with $\Phi \in \mathcal{A}$. Assume that there exists an $(L, \|\cdot\|, |\cdot|)$ -Lipschitz function, h , such that $h(x) \geq 0$ if and only if $\Phi(x) = \text{True}$. Then for ψ , there also exists an $(L, \|\cdot\|^{t_f}, |\cdot|)$ -Lipschitz robustness measure, $\rho : \mathcal{S}_{t_f}^{\mathbb{R}^n} \rightarrow \mathbb{R}$, as in (11), i.e.,*

$$|\rho(s') - \rho(s)| \leq L \max_{0 \leq t \leq t_f} \|s'(t) - s(t)\| = L \|s' - s\|^{t_f}. \quad (32)$$

Proof: It is known from prior works, e.g. [17], that the following are valid robustness measures for specifications, ψ , satisfying Assumption 1, as $h(x) \geq 0 \iff \Phi(x) = \text{True}$:

$$\rho(s) = \begin{cases} \max_{0 \leq t \leq t_f} h(s(t)), & \text{If } \psi = \text{True} \cup \Phi, \\ \min_{0 \leq t \leq t_f} h(s(t)), & \text{If } \psi = \neg(\text{True} \cup \neg\Phi). \end{cases}$$

We will show that the first measure is Lipschitz.

$$|\rho(s') - \rho(s)| = \left| \max_{0 \leq t \leq t_f} h(s'(t)) - \max_{0 \leq t \leq t_f} h(s(t)) \right|, \quad (33)$$

$$\leq \max_{0 \leq t \leq t_f} |h(s'(t)) - h(s(t))|, \quad (34)$$

$$\leq L \max_{0 \leq t \leq t_f} \|s'(t) - s(t)\|. \quad (35)$$

Following a similar chain of logic shows that the second measure is Lipschitz, thus completing the proof. ■

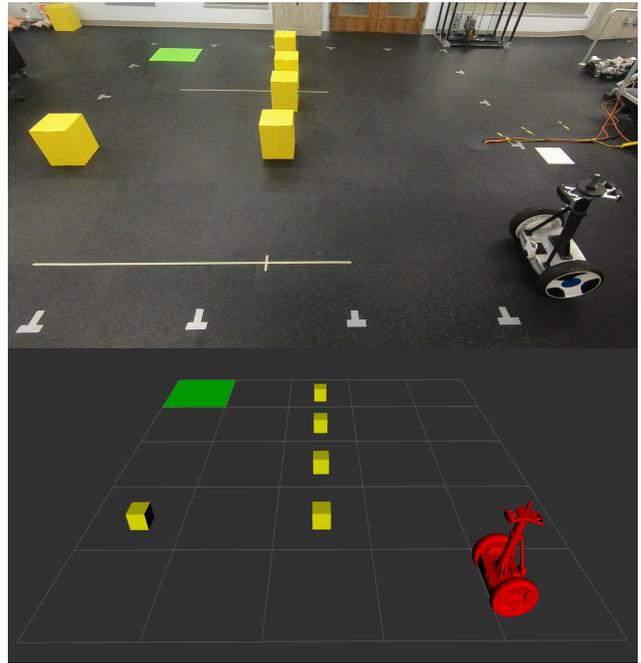


Fig. 2. A picture of the (bottom) ROS-based simulation environment we used to approximate our (top) true, Segway system.

Finally, Propositions 1 and 2 ensure that there always exists a Lipschitz continuous robustness measure, ρ , for any specification, ψ , satisfying the conditions for Assumption 1.

Proposition 3. *For any specification, ψ , that satisfies Assumption 1, there always exists a Lipschitz continuous robustness measure, ρ , of the form in (11).*

Proof: The proof is a direct application of Propositions 1 and 2. Proposition 1 guarantees the existence of at least one $(L, \|\cdot\|, |\cdot|)$ -Lipschitz function, h , for $\Phi \in \mathcal{A}$, $L \in \mathbb{R}_+$, and norm, $\|\cdot\|$, on \mathbb{R}^n . Using this function, h , in the associated robustness measure, ρ , for ψ , as specified in Proposition 2, guarantees the existence of at least one Lipschitz robustness measure. ■

As a result, Proposition 3 indicates that, for any specification, ψ , that satisfies Assumption 1, there always exists a Lipschitz continuous robustness measure, ρ . As a result, if the nominal system is $(\epsilon, t_f, \lambda, \|\cdot\|)$ -accurate, with respect to the same norm, $\|\cdot\|$, that induces the norm for which ρ is Lipschitz, $\|\cdot\|^{t_f}$, the results of Theorem 1 still hold.

IV. EXPERIMENTAL RESULTS

Our goal is to verify that a Nonlinear MPC controller, designed to steer a Segway to a goal while avoiding obstacles, satisfies a safety specification, ψ , incumbent on the Segway. This verification question arose when, under normal operation, this controller appeared to successfully steer the Segway to the goal while also satisfying ψ . Hence, we want to use Theorem 1, to verify that this controller always satisfies ψ , regardless of the goals provided to the Segway. The experimental setup and accompanying simulation, for this verification procedure, is shown in Figure 2.

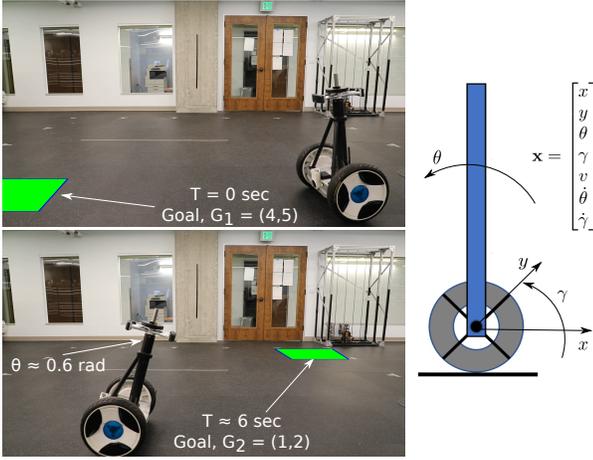


Fig. 3. Shown above is the Segway operating with an example environmental configuration, $d = [4, 5, 1, 2, 4.885]^T$. Each vector, $d = [G_1, G_2, T]^T \in \mathbb{R}^5$ denotes the first goal cell to which the Segway is to navigate, $G_1 \in \{1, 2, 3, 4, 5\}^2$, the second goal cell, $G_2 \in \{1, 2, 3, 4, 5\}^2$, and the time, $T \in [0, 10]$, the Segway is to switch between going to G_1 and going to G_2 . The example shown is a worst-case configuration, the experimental trajectory for which is shown, in red, in Figure 6.

To use Theorem 1 for this verification, we require a feasible space of environmental configurations, \mathcal{D} . As we are trying to verify the controller's capacity to satisfy ψ regardless of the goals provided to it, we specify a vector, $d \in \mathcal{D} \subset \mathbb{R}^5$, to be a vector of two goals, G_1, G_2 , and a switching time, T . This time, T , determines when the Segway switches from navigating to the first goal, G_1 , to navigating to the second goal, G_2 . Goals are represented as shaded, cellular regions, as in the green region in Figure 2. Hence, $G_1, G_2 \in \{1, 2, 3, 4, 5\}^2$ are cells in this 5×5 grid, and the switching time, $T \in [0, 10]$. Figure 3 portrays a graphical representation of the Segway undergoing a specific environmental configuration, and the full, mathematical setup is:

$$x = [x, y, \theta, \psi, v, \dot{\theta}, \dot{\psi}] \in \mathbb{R}^7, \quad d = [G_1, G_2, T]^T \in \mathcal{D} \subset \mathbb{R}^5, \\ \llbracket \Phi \rrbracket = \{\theta \in \mathbb{R} \mid |\theta| \leq 0.7 \text{ rad}\}, \quad \psi = \neg(\text{True} \cup \neg\Phi).$$

Given the above environment and safety specification, ψ , the measurement function h and the robustness measure are defined as:

$$h(x) = 0.7 - |\theta|, \quad (36)$$

$$\rho(\bar{\phi}) = \min_{0 \leq t \leq t_f} h(\bar{\phi}(t)). \quad (37)$$

Here, we abbreviate $\hat{\phi}(x_0, U(x(t), d), d) = \bar{\phi}$. With this setup, Theorem 1 requires the parameters, ϵ, t_f , and λ , for our ROS-based simulator, and the nominal worst-case robustness h^* as in (20), for the robustness measure (37).

To determine the parameters for our simulator, we note that our measurement function (36) is $(1, \|\cdot\|_\alpha, |\cdot|)$ -Lipschitz with respect to the norm,

$$\|x\|_\alpha = \sum_{i=1}^7 \alpha_i x_i^2, \quad \alpha_i = \begin{cases} 1 & \text{If } i = 3 \\ 10^{-6} & \text{Else.} \end{cases}$$

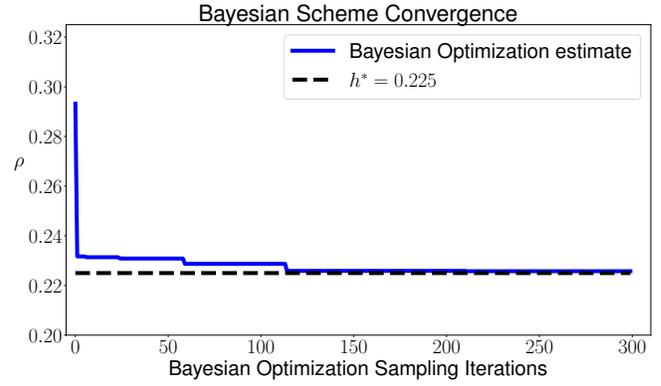


Fig. 4. Shown above is the convergence of the Bayesian-Testing framework, minimizing the robustness measure, ρ , stated in Section IV, over the ROS-based simulation of our Segway. Notice that by 125 iterations, the Bayesian optimization estimate converged to $h^* = 0.225$.

Hence, robustness measure (37) is also $(1, \|\cdot\|_\alpha^t, |\cdot|)$ -Lipschitz. As Theorem 1 requires that our simulator be accurate with respect to the same norm, $\|\cdot\|_\alpha$, that induces the norm over which ρ is Lipschitz, $\|\cdot\|_\alpha^t$, we determined the accuracy of our simulator with respect to this norm, $\|\cdot\|_\alpha$. To do so, we recorded the maximum simulation error, $\Delta = \max_{0 \leq t \leq t_f} \|\Delta_\phi(x_0, U, d, w(t), t)\|_\alpha$ (17), for $N = 300$ different pairs of simulation and system trajectories, wherein the starting, goal, and obstacle locations were the same for each pair. The recorded data is shown in Figure 5. Here, we inherently assume that the maximum norm deviance, Δ , is a random variable distributed by some unknown distribution, π_Δ . As we do not have access to π_Δ , we construct a Monte-Carlo based estimate probability distribution, π_Δ^N . This distribution, π_Δ^N , is known to be an unbiased estimator of π_Δ , for which the absolute value of the difference in variance of the two probability measures is bounded above by $\frac{1}{N}$. Based on this estimate distribution, π_Δ^N , our data indicates that our ROS-based simulator is $(\epsilon = 0.125, t_f = 10, \lambda = 0.05, \|\cdot\|_\alpha)$ -accurate, as $\mathbb{P}[\Delta \leq \epsilon = 0.125] \geq (1 - \lambda = 0.95)$.

It remains to determine, for this simulator, the nominal worst-case robustness measure, h^* , defined in equation (20), with respect to the robustness measure (37). To solve optimization problem (20), we used the Bayesian optimization scheme detailed in [22], which was proven to converge in expectation. Figure 4 shows the convergence results of this algorithm, when run for 300 iterations, in attempting to solve optimization problem (20). By roughly 125 iterations, the algorithm had converged to an $h^* = 0.225$, which, based on the convergence results of this algorithm in [22], we assume to be the solution to optimization problem (20). Furthermore, this h^* arose for the environmental configuration, $d = [G_1 = (4, 5), G_2 = (1, 2), T = 4.885]^T$.

As a result, Theorem 1 indicates that, as $(h^* = 0.225) \geq (L\epsilon = 0.125)$,

$$\mathbb{P}_w [\tilde{\phi}(t_f) \models \psi] \geq 0.95 \quad \forall d \in \mathcal{D}. \quad (38)$$

Here, we abbreviated $\phi(x_0, U(x(t), d), d, w(t)) = \tilde{\phi}$. To check the accuracy of this probabilistic claim, seven, independent runs of the Segway undergoing this worst-case

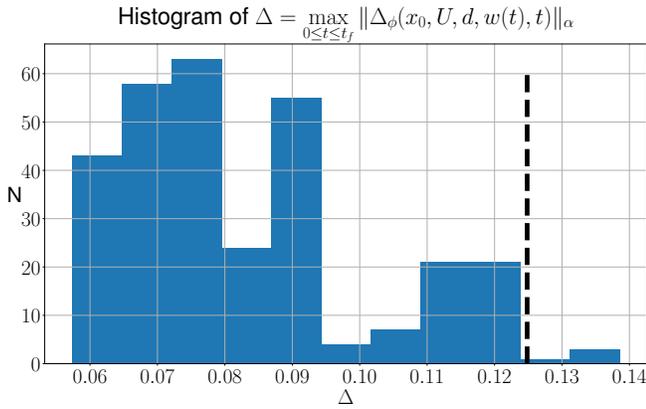


Fig. 5. Shown above is a histogram of the sampled, upper-bound on angular trajectory deviation, Δ , with respect to the norm, $\|\cdot\|_\alpha$. The dashed, black line indicates the cutoff ϵ value, $\epsilon = 0.125$, for which $\mathbb{P}[\Delta \leq \epsilon] \geq 0.95$.

environment configuration, $d = [4, 5, 1, 2, 4.885]$, are shown in Figure 6. The probabilistic verification claim (38) indicates that, with probability at least 0.95, experimental trajectories should satisfy ψ , as they should all lie within the lavender region - and they indeed all lie within that box. For a video summarizing this paper, in addition to experimental demonstrations, please reference [1].

V. CONCLUSION

Our main contribution extends prior work done in the Bayesian testing community, by describing a method to lower bound the probability by which a system will satisfy an operational specification in practice. Our method constructs a probabilistic signal tube around a simulation trajectory, in which we anticipate the real system trajectory to lie, with some minimum probability. Then, we determine the worst-case robustness measure for this nominal system, which, in conjunction with the signal tube, we use to lower bound the worst case robustness measure for the true system. If this lower bound is positive, the controller is verified to satisfy its specification, with the same minimum probability.

REFERENCES

- [1] "Video of experiment." <https://youtu.be/qwc7Jei-FxE>.
- [2] P. Nilsson, O. Hussien, Y. Chen, A. Balkan, M. Rungger, A. Ames, J. Grizzle, N. Ozay, H. Peng, and P. Tabuada, "Preliminary results on correct-by-construction control software synthesis for adaptive cruise control," in *53rd IEEE Conference on Decision and Control*, pp. 816–823, IEEE, 2014.
- [3] P. Nilsson, O. Hussien, A. Balkan, Y. Chen, A. D. Ames, J. W. Grizzle, N. Ozay, H. Peng, and P. Tabuada, "Correct-by-construction adaptive cruise control: Two approaches," *IEEE Transactions on Control Systems Technology*, vol. 24, no. 4, pp. 1294–1307, 2015.
- [4] S. A. Seshia, D. Sadigh, and S. S. Sastry, "Towards verified artificial intelligence," *arXiv preprint arXiv:1606.08514*, 2016.
- [5] S. Prajna and A. Jadbabaic, "Safety verification of hybrid systems using barrier certificates," in *International Workshop on Hybrid Systems: Computation and Control*, vol. 2993, pp. 477–492, Springer Verlag, 2004.
- [6] J. Kapinski, S. Sankaranarayanan, J. V. Deshmukh, and N. Aréchiga, "Simulation-guided Lyapunov analysis for hybrid dynamical systems," in *HSCC 2014 - Proceedings of the 17th International Conference on Hybrid Systems: Computation and Control (Part of CPS Week)*, pp. 133–142, Association for Computing Machinery, 2014.

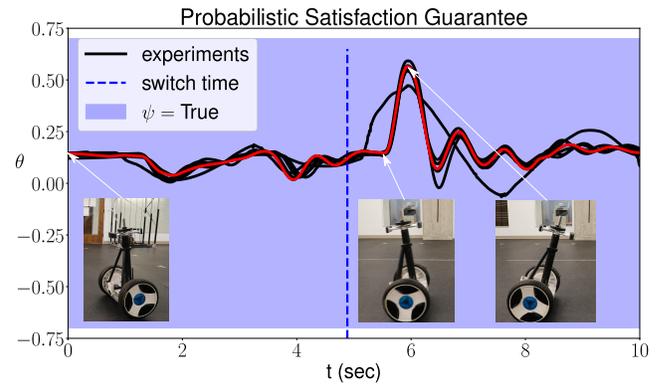


Fig. 6. Shown above is the angular trajectories of the Segway undergoing seven independent trials of the worst case environmental configuration derived in Section IV. The goal switching time, T , is demarcated by the dashed, blue line. Several snapshots of the Segway, during the trajectory highlighted in red, are shown at the bottom. We anticipate the system to lie within the lavender region shown, with probability ≥ 0.95 , and as such, satisfy its safety specification, ψ , to the same minimum probability.

- [7] 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*, vol. 62, no. 8, pp. 3861–3876, 2016.
- [8] M. Althoff and S. Lutz, "Automatic generation of safety-critical test scenarios for collision avoidance of road vehicles," in *2018 IEEE Intelligent Vehicles Symposium (IV)*, pp. 1326–1333, IEEE, 2018.
- [9] M. Koschi, C. Pek, S. Maierhofer, and M. Althoff, "Computationally efficient safety falsification of adaptive cruise control systems," in *2019 IEEE Intelligent Transportation Systems Conference (ITSC)*, pp. 2879–2886, IEEE, 2019.
- [10] T. A. Wheeler and M. J. Kochenderfer, "Critical factor graph situation clusters for accelerated automotive safety validation," in *2019 IEEE Intelligent Vehicles Symposium (IV)*, pp. 2133–2139, IEEE, 2019.
- [11] N. Hansen, "The cma evolution strategy: A tutorial," *arXiv preprint arXiv:1604.00772*, 2016.
- [12] Y. Annappureddy, C. Liu, G. Fainekos, and S. Sankaranarayanan, "S-taliro: A tool for temporal logic falsification for hybrid systems," in *International Conference on Tools and Algorithms for the Construction and Analysis of Systems*, pp. 254–257, Springer, 2011.
- [13] R. Calandra, A. Seyfarth, J. Peters, and M. P. Deisenroth, "Bayesian optimization for learning gaits under uncertainty," *Annals of Mathematics and Artificial Intelligence*, vol. 76, no. 1, pp. 5–23, 2016.
- [14] F. Berkenkamp, A. Krause, and A. P. Schoellig, "Bayesian optimization with safety constraints: safe and automatic parameter tuning in robotics," *arXiv preprint arXiv:1602.04450*, 2016.
- [15] S. Ghosh, F. Berkenkamp, G. Ranade, S. Qadeer, and A. Kapoor, "Verifying controllers against adversarial examples with bayesian optimization," in *2018 IEEE International Conference on Robotics and Automation (ICRA)*, pp. 7306–7313, IEEE, 2018.
- [16] B. Gangopadhyay, S. Khastgir, S. Dey, P. Dasgupta, G. Montana, and P. Jennings, "Identification of test cases for automated driving systems using bayesian optimization," in *2019 IEEE Intelligent Transportation Systems Conference (ITSC)*, pp. 1961–1967, IEEE, 2019.
- [17] C. Madsen, P. Vaidyanathan, S. Sadraddini, C.-I. Vasile, N. A. DeLateur, R. Weiss, D. Densmore, and C. Belta, "Metrics for signal temporal logic formulae," in *2018 IEEE Conference on Decision and Control (CDC)*, pp. 1542–1547, IEEE, 2018.
- [18] A. Rizk, G. Batt, F. Fages, and S. Soliman, "Continuous valuations of temporal logic specifications with applications to parameter optimization and robustness measures," *Theoretical Computer Science*, vol. 412, no. 26, pp. 2827–2839, 2011.
- [19] Y. Gilpin, V. Kurtz, and H. Lin, "A smooth robustness measure of signal temporal logic for symbolic control," *arXiv preprint arXiv:2006.05239*, 2020.
- [20] C. Baier and J.-P. Katoen, *Principles of model checking*. MIT press, 2008.
- [21] P. I. Frazier, "A tutorial on bayesian optimization," *arXiv preprint arXiv:1807.02811*, 2018.
- [22] A. D. Bull, "Convergence rates of efficient global optimization algorithms," *Journal of Machine Learning Research*, vol. 12, no. 10, 2011.