

Requirements-driven Test Generation for Autonomous Vehicles with Machine Learning Components*

Cumhur Erkan Tuncali¹, Georgios Fainekos¹, Danil Prokhorov², Hisahiro Ito², and James Kapinski²

Abstract—Autonomous vehicles are complex systems that are challenging to test and debug. A requirements-driven approach to the development process can decrease the resources required to design and test these systems, while simultaneously increasing the reliability. We present a testing framework that uses signal temporal logic (STL), which is a precise and unambiguous requirements language. Our framework evaluates test cases against the STL formulae and additionally uses the requirements to automatically discover test cases that fail to satisfy the requirements. One of the key features of our tool is the support for machine learning (ML) components in the system design, such as deep neural networks. The framework allows evaluation of the control algorithms, including the ML components, and it also includes models of CCD camera, lidar, and radar sensors, as well as the vehicle environment. We use multiple methods to generate test cases, including covering arrays, which is an efficient method to search discrete variable spaces. The resulting test cases can be used to debug the controller design by identifying controller behaviors that do not satisfy requirements. The test cases can also enhance the testing phase of development by identifying critical corner cases that correspond to the limits of the system’s allowed behaviors. We present STL requirements for an autonomous vehicle system, which capture both component-level and system-level behaviors. Additionally, we present three driving scenarios and demonstrate how our requirements-driven testing framework can be used to identify critical system behaviors, which can be used to support the development process.

I. INTRODUCTION

Autonomous driving systems are in a stage of rapid research and development spanning a broad range of maturity from simulations to on-road testing and deployment. They are expected to have a significant impact on the vehicle market and the broader economy and society in the future.

Testing of highly automated and autonomous driving systems is also an area of active research. Both governmental and non-governmental organizations are grappling with the unique requirements of these new, highly complex systems, as they have to operate safely and reliably in diverse driving environments. Government and industry sponsored partnerships have produced a number of guiding documents and clarifications, such as NHTSA [1], SAE [2], CAMP [3], NCAP [4], PEGASUS [5]. The research community has also been contributing to the development of methodologies for testing automated driving systems.

Stellet et al. [6] surveyed existing approaches to testing such as simulation-only, X-in-the-loop and augmented reality approaches, as well as test criteria and metrics (see also [7]). Koopman and Wagner identified challenges of testing

and proposed potential solutions, such as fault injection, as a way to perform more efficient edge case testing [8]. The publications [9] and [10] provide in-depth discussions on the challenges of safety validation for autonomous vehicles, arguing that virtual testing should be the main target for both methodological and economic reasons.

However, no universally agreed upon testing or verification methods have yet arisen for autonomous driving systems. One reason is that the current autonomous systems architectures usually include some Machine Learning (ML) components, such as Deep Neural Networks (DNNs), which are notoriously difficult to test and verify. For instance, Automated Driving System (ADS) designs often use ML components such as DNNs to classify objects within CCD images and to determine their positions relative to the vehicle, a process known as *object detection and classification* [11], [12]. Other designs use Neural Networks (NNs) to perform *end-to-end* control of the vehicle, meaning that the NN takes in the image data and outputs actuator commands, without explicitly performing an intermediate object detection step [13], [14], [15]. Still other approaches use end-to-end learning to do intermediate decisions like risk assessment [16].

ML system components are problematic from an analysis perspective, as it is difficult or impossible to characterize all of the behaviors of these components under all circumstances. One reason is that the complexity of these systems can be very high in terms of the number of parameters. For example, AlexNet [17], a pre-trained DNN that is used for classification of CCD images, has 60 million parameters. Another reason for the difficulty in characterizing behaviors of ML components is that the parameters are learned based on training data. In other words, characterizing ML behaviors is, in some ways, as difficult as the task of characterizing the training data. Again using the AlexNet example, the number of training images used was 1.2 million. While the main strength of DNNs is their ability to generalize from training data, the major challenge for analysis is that we do not understand well how they generalize to all possible cases. Therefore, there has been significant interest on verification and testing for ML components. For example, adversarial testing approaches seek to identify perturbations in image data that result in misclassifications [18], [19], [20]. However, most of the existing work on testing and verification of systems with ML components focuses only on the ML components themselves, without consideration of the closed-loop behavior of the system.

The closed-loop nature of a typical autonomous driving system can be described as follows. A *perception* system processes data gathered from various sensing devices, such as cameras, lidar, and radar. The output of the perception system is an estimation of the principal (*ego*) vehicle’s position

*This work was partially funded by NSF awards CNS 1446730, 1350420

¹School of Computing, Informatics & Decision Systems Engineering, Arizona State University, USA etuncali, fainekos@asu.edu

²Toyota Technical Center, Ann Arbor MI, USA danil.prokhorov@toyota.com

with respect to external obstacles (e.g., other vehicles, called *agent* vehicles, and pedestrians). A *path planning* algorithm uses the output of the perception system to produce a short-term plan for how the ego vehicle should behave. A *tracking controller* then takes the output of the path planner and produces actuation outputs, such as accelerator, braking, and steering commands. The actuation commands affect the vehicle’s interaction with the environment. The iterative process of sensing, processing, and actuating is what we refer to as closed-loop behavior.

It is important to note that by design, closed-loop systems have error tolerance mechanisms. Hence, adversarial attacks that work on individual ML components may not have the same effect on the closed loop system. Since for ADS applications the ultimate goal is to evaluate the closed-loop system performance, any testing methods used to evaluate such systems should support this goal.

We present a framework for Simulation-based Adversarial Testing of Autonomous Vehicles (Sim-ATAV), which can be used to check closed-loop properties of ADS that include ML components. In particular, our work focuses on methods to determine perturbations in the configuration of a test scenario, meaning that we seek to find scenarios that lead to unexpected behaviors, such as misclassifications and ultimately vehicle collisions. The framework that we present allows this type of testing in a virtual environment. By utilizing advanced 3D models and image rendering tools, such as the ones used in game engines, the gap between testing in a virtual environment and the real world can be minimized. We describe a testing methodology, based on a test case generation method, called covering arrays [21], and requirement falsification methods [22] to automatically identify problematic test scenarios. The resulting framework can be used to increase the reliability of autonomous driving systems.

An earlier version of this work appeared in [23]. The contributions of that work can be summarized as follows. In [23], we provided a new algorithm to perform falsification of formal requirements for an autonomous vehicle in a closed-loop with the perception system, which includes an efficient means of searching over discrete and continuous parameter spaces. The method represents a new way to do adversarial testing in scenario *configuration space*, as opposed to the usual method, which considers adversaries in *image space*. Additionally, we demonstrated a new way to characterize problems with perception systems in *configuration space*. Lastly, we extended the software testing theory of covering arrays to closed-loop Cyber-Physical System (CPS) applications that have embedded ML algorithms.

The present paper provides the following contributions that are in addition to those from [23]:

- We add models of lidar and radar sensors and include sensor fusion algorithms, and we demonstrate how the requirements-based testing framework we propose can be used to automate the search for specific types of fault cases involving sensor interactions.
- We provide requirements for both component-level and system-level behaviors, and we show how to automate the identification of behaviors where component-level

failures lead to system-level failures. An example of the kind of analysis this allows is automatically finding cases where a sensor failure leads to a collision case.

- We include a model of agent *visibility* to various sensors and include this notion in the requirements that we consider. This provides a way to reason about how the system should behave, based on whether agents are or are not visible, including the ability to reason about the temporal aspects of agent visibility. For example, we can use this feature to test the requirement that within 1 second after an agent becomes visible to the lidar sensor, the perception system should correctly classify the agent. This allows us to automate the search for behaviors related to temporal aspects of sensor behaviors in the context of a realistic driving scenario.
- We demonstrate the ability to falsify properties by adversarially searching over agent trajectories. This permits the use of our requirements-driven search-based approach over a broad class of agent behaviors, which allows us to automatically identify corner cases that are difficult to find using traditional simulation-based techniques.
- We have released a publicly available toolbox Sim-ATAV as an add-on to the S-TaLiRo falsification toolbox for Matlab[®] [24]: <https://sites.google.com/asu.edu/s-taliro/s-taliro/sim-atav>

II. RELATED WORK

Testing and evaluation methods for Autonomous¹ Vehicles (AVs) could be categorized into three major classes: (1) model based, (2) data-driven, and (3) scenario based. Scenario-based approaches utilize accident reports and driving conditions that are easily identifiable as challenging, producing specific test scenarios to be executed either in the real world or in a simulation environment. For example, Euro NCAP [4] and DOT [25] provide such scenarios. Data-driven approaches, on the other hand, typically utilize driving data [26] to generate probabilistic models of human drivers. Such models are then used for risk assessment and rare event sampling for AV algorithms under specific driving scenarios [27].

The aforementioned testing methods are important and necessary before AV deployment, but they cannot help with design exploration and automated fault detection at early development stages. Such problems are addressed by model-based verification [28], [29], model based test generation [30], [31], [32], [33], [34], [35], [36], or a combination thereof [37], [38]. It is important to also highlight that these methods typically ignore or use simple models to abstract away proximity sensors and, especially, the vision systems. However, ignoring sensors or using simplified sensing models may be a dangerously simplifying assumption since it ignores the complex interactions between the dynamics of the vehicle and the sensors. For example, the effective sensing range of a sensor platform mounted on the roof of a vehicle is affected when the vehicle makes hard turns.

¹We utilize the more general term “autonomous” as opposed to a more restricted “automated” since our methods could potentially apply to all levels of autonomy.

In addition, vision-based perception systems have become an integral component of the sensor platform of AVs, and in many cases, they constitute the only perception system. Currently, the winning algorithmic technology for image processing systems is utilizing DNNs. For instance, by 2011, the DNN architecture proposed in [39] was already capable of classifying pre-segmented images of traffic signs with better accuracy than humans (99.46% vs 99.22%). Since then, there has been substantial progress with DNNs performing both segmentation and classification [40], [41]. Yet, in spite of the multiple impressive results using DNNs, it is still also easy to devise methods that can produce (so-called adversarial) images that will fool them [18], [19].

The latter (negative) result raises two important questions: (1) can we still generate adversarial inputs for DNNs when we manipulate the physical properties and trajectories of the objects in the environment of the AV, and (2) how does the DNN accuracy affect the system level properties of an AV, that is, its functional safety? Exhaustive verification methods for DNNs in-the-loop are still in their infancy [42], and they cannot handle AVs with DNN components in the loop. To address the two questions above, several model-based test generation methods have been proposed [43], [44], [45], [23]. The procedure described in [43], [44] analyzes the performance of the perception system using static images to identify candidate counterexamples, which are then checked using simulations of the closed-loop system to determine whether the AV exhibits unsafe behaviors. On the other hand, [45], [23] develop methods that directly search for unsafe behaviors of the closed-loop system by defining a cost function on the closed-loop behaviors. The differences between [45] and [23] are primarily on the search methods, the simulation environments, and the AVs, with [23] providing a more efficient method for combinatorial search.

In this extended version of [23], we take the system-level adversarial test generation methods for AV one step further. We demonstrate that our framework [23] can be extended for test generation for AV with multi-sensor systems as opposed to vision-only perception systems. Moreover, we demonstrate the importance and effectiveness of test generation methods guided by system-level requirements as well as perception-level requirements.

Using our framework, we can formalize and test against requirements on the sensor performance in the context of a driving scenario. For example, the lidar's point cloud density drops significantly with the distance to the target object, for example, a pedestrian. Similar to this aspect of lidar behavior, the pixel count of a CCD camera would also decrease dramatically with the distance if it were to be used for pedestrian detection, since the area of an observed object decreases as the square of the distance to the object. This may complicate testing for long-range observation conditions. Our framework supports testing these aspects of sensor performance.

III. PRELIMINARIES

This section presents the setting used to describe the testing procedures performed with our framework. The purpose of

our framework is to provide a mechanism to test, evaluate, and improve on an autonomous driving system design. To do this, we use a simulation environment that incorporates models of a vehicle (called the *ego* vehicle), a perception system, which is used to estimate the state of the vehicle with respect to other objects in its environment, a controller, which makes decisions about how the vehicle will behave, and the environment in which the ego vehicle is deployed. The environment model contains representations of a wide variety of objects that can interact with the ego vehicle, including roads, buildings, pedestrians, and other vehicles (called *agent* vehicles). The behaviors of the system are determined by the evolution of the model states over time, which we compute using a *simulator*.

In the following, \mathbb{R} represents the set of real numbers, while \mathbb{Z} the set of integers. In addition, $\bar{\mathbb{R}} = \mathbb{R} \cup \{\pm\infty\}$ and $\mathbb{R}_{\geq 0}$ is the set of positive reals. Formally, we assume that a test scenario Σ is captured by a simulation function $sim : \mathcal{X}_0 \times \mathbf{U} \times \mathcal{P} \times \mathbb{R}_{\geq 0} \rightarrow \mathbf{Y}$ that maps a vector of initial conditions $x_0 \in \mathcal{X}_0$, a vector of parameters $p \in \mathcal{P}$, a total simulation time $T \in \mathbb{R}_{\geq 0}$ and a time stamped input signal $\mathbf{u} \in \mathbf{U}$ to a time stamped output signal $\mathbf{y} \in \mathbf{Y}$.

Here, $\mathcal{X}_0 \subseteq \mathbb{R}^{n_1^x} \times \mathbb{Z}^{n_2^x}$ is the set of initial conditions for the whole scenario, i.e., for the ego vehicle(s) as well as any other stateful object in the environment. The variable n_1^x captures the number of the continuous-valued state variables in the system (i.e., the order of the differential and/or difference equations), and n_2^x captures the number of discrete-valued (and, primarily, finite-valued) state variables in the system. In other words, we assume that the models we consider are hybrid dynamical systems [46]. Similarly, $\mathcal{P} \subseteq \mathbb{R}^{n_1^p} \times \mathbb{Z}^{n_2^p}$ is a set of n_1^p continuous-valued parameters, such as ambient temperature, light intensity, or color, and n_2^p discrete-valued (categorical) parameters, such as vehicle model or sign type.

The set of potential input values is denoted by $\mathcal{U} \subseteq \mathbb{R}^{m_{in}}$, where m_{in} is the number of time varying input signals to the test scenario. The set of possible input signals is $\mathbf{U} = (\mathcal{U} \times \mathbb{R}_{\geq 0})^{N_{in}+1}$, where N_{in} is the number of samples for the input signal. In other words, an *input signal* (also referred to as *input trace*) $\mathbf{u} \in \mathbf{U}$ is a function $\mathbf{u} : \{0, 1, \dots, N_{in}\} \rightarrow \mathcal{U} \times \mathbb{R}_{\geq 0}$ which maps each sample i to an input value $u_i \in \mathcal{U}$ and a time stamp $t_i \in \mathbb{R}_{\geq 0}$. Alternatively, we can view \mathbf{u} as a finite sequence of input signal values and their corresponding times:

$$\mathbf{u} = (u_0, t_0)(u_1, t_1) \cdots (u_{N_{in}}, t_{N_{in}}).$$

Here, we make two assumptions : (i) $t_{N_{in}}$ is no greater than the simulation time T (i.e., $t_{N_{in}} \leq T$), and (ii) the timestamps are monotonically increasing: $\forall i, j \in \{0, 1, \dots, N_{in}\}$, if $i < j$, then $t_i < t_j$. Finally, since the simulator may need to produce output values at some time t between two timestamps (i.e., $t \in (t_i, t_{i+1})$), we will assume that the simulator decides what interpolation function it will use.

Given initial conditions $x_0 \in \mathcal{X}_0$, parameter values $p \in \mathcal{P}$, input signals $\mathbf{u} \in \mathbf{U}$, and the total simulation time T , the simulator returns an *output trajectory* (also referred to as *output trace*) $\mathbf{y} = sim(x_0, \mathbf{u}, p, T)$. The output trace is a function $\mathbf{y} : \{0, 1, \dots, N_{out}\} \rightarrow \mathcal{Y} \times \mathbb{R}_{\geq 0}$ that maps each sample i to an output value $y_i \in \mathcal{Y} \subseteq \mathbb{R}^{m_{out}}$ and a time



Fig. 1: The simple test scenario of Example 1.

stamp $t_i \in \mathbb{R}_{\geq 0}$. Here, m_{out} is the number of observable output variables. We denote the set of possible output traces by $\mathbf{Y} = (\mathcal{Y} \times \mathbb{R}_{\geq 0})^{N_{out}}$. The output trace timestamps should satisfy (i) $t_{N_{out}} = T$, and (ii) the monotonicity property.

For notational convenience, we will make an additional assumption that the simulator also returns an updated input trace \mathbf{u} where $N = N_{in} = N_{out}$ and the timestamps of \mathbf{y} and \mathbf{u} match. We refer to the triple $\sigma = (\mathbf{y}, \mathbf{u}, p)$ as a *simulation trace*, which can also be viewed as a function $\sigma : \{0, 1, \dots, N_{out}\} \rightarrow \mathcal{Y} \times \mathcal{U} \times \mathcal{P} \times \mathbb{R}_{\geq 0}$, or as a sequence (recall that p is constant):

$$\sigma = (y_0, u_0, p, t_0)(y_1, u_1, p, t_1) \cdots (y_N, u_N, p, t_N).$$

We denote the set of all simulation traces σ of Σ by $\mathcal{L}(\Sigma)$.

Remark 1: In this paper, the function *sim* is assumed to be deterministic; however, the results we present are also applicable to stochastic systems (i.e., when the *sim* function is stochastic). See [47] for a discussion.

Example 1: We will present a simple illustrative example to clarify the notation. Let's assume a test scenario as in Fig. 1. For $i \in \{a, e\}$, we will denote by $z^{(i)}$ the longitudinal position of the vehicle i and by $v^{(i)}$ the velocity of the vehicle i . The ego vehicle (e) implements an adaptive cruise control (ACC) algorithm, which we treat as a black box: $\dot{\eta}^{(e)} = f_e(\eta^{(e)}, \eta^{(a)})$, where $\eta^{(i)} = [z^{(i)} v^{(i)}]^T$. In this simple model of an ego car, the ego car senses its environment by measuring the state $\eta^{(a)}$ of the adversarial agent. The adversarial agent (a) has simple integrator dynamics $\dot{z}^{(a)} = \mu v^{(a)}$ and $\dot{v}^{(a)} = \xi$ (with the additional constraint of no negative velocity, i.e., $v(t) \geq 0$), where ξ and μ are the time varying inputs that we search over, i.e., $u = [\xi \mu]^T$. In particular, $\mu(t) \in \{1, 2\}$ models the normal versus the sport driving mode in the powertrain (selected by the driver), while $\xi(t) = [-1, 1]$ is the acceleration (and braking) input also provided by the driver of the adversarial vehicle. Assuming $N_{in} = 201$, then $\mathbf{U} = ([-1, 1] \times \{1, 2\})^{201}$ – see Fig. 2 for an example input. In this test scenario, the state space is $x = [z^{(e)} v^{(e)} z^{(a)} v^{(a)}]^T$ and, thus \mathcal{X}_0 is the set of initial positions $z^{(i)}$ and velocities $v^{(i)}$ of the two vehicles. The set of parameters is empty since this test scenario does not have any constant parameters. The output trace is defined to be the positions of the two vehicles over time, i.e., $y = [z^{(e)} z^{(a)}]^T$. Fig. 3 presents the vehicle positions over time for the inputs in Fig. 2. It can be observed that the ego vehicle does not utilize a safe ACC since it collides with the adversarial vehicle at about time 4.

A. Signal Temporal Logic

Signal Temporal Logic (STL) was introduced as a syntactic extension to Metric Temporal Logic (MTL) to reason about

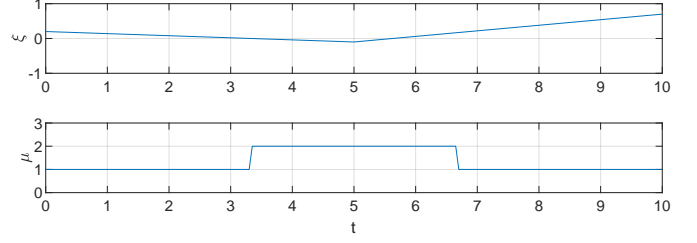


Fig. 2: Input trace u for the test scenario of Example 1.

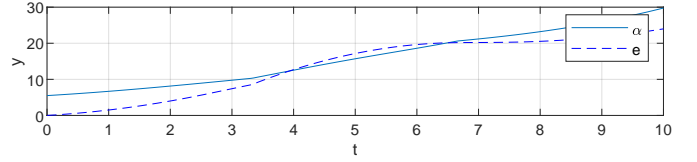


Fig. 3: Output trace y for the test scenario of Example 1.

real-time properties of signals (simulation traces) (for an overview see [48]). STL formulae are built over predicates on the variables of a signal using Boolean and temporal operators. The temporal operators include *eventually* ($\diamond_{\mathcal{I}}$), *always* ($\square_{\mathcal{I}}$) and *until* ($U_{\mathcal{I}}$), where \mathcal{I} is a time interval that encodes timing constraints. The boolean operators include *conjunction* \wedge , *disjunction* \vee , *negation* \neg , and *implication* \implies .

In this work, we interpret STL formulas over the observable simulation traces. STL specifications can describe the usual properties of interest in system design such as (bounded time) **reachability**, for example, *eventually, between time 1.2 and 5 (not including), y should drop below -10* : $\diamond_{[1.2, 5)}(y \leq -10)$, and **safety**, for example, *after time 2 time units, y should always be greater than 10*: $\square_{[2, +\infty)}(y \geq 10)$. An important class of expressible requirements in STL are **reactive requirements**, such as $\square((y \leq -10) \rightarrow \diamond_{[0, 2]}(y \geq 10))$, which states that *whenever y drops below -10 , then within 2 time units y should rise above 10*.

Informally speaking, we allow predicate expressions to capture arbitrary constraints over the output variables, inputs, and parameters of the system. More formally, we assume that predicates π are expressions built using the grammar $\pi ::= f(y, u, p) \geq c \mid \neg \pi_1 \mid (\pi) \mid \pi_1 \vee \pi_2 \mid \pi_1 \wedge \pi_2$, where f is a function and c is a constant in \mathbb{R} . In other words, each predicate π represents a subset in the space $\mathcal{Y} \times \mathcal{U} \times \mathcal{P}$. In the following, we represent the set that corresponds to the predicate π using the notation $\mathcal{O}(\pi)$. For example, if $\pi = (y^{(1)} \leq -10) \vee (y^{(1)} + y^{(2)} \geq 10)$ where $y^{(i)}$ is the i -th component of the vector y , then $\mathcal{O}(\pi) = (\infty, -10] \times \mathbb{R} \cup \{y \in \mathbb{R}^2 \mid y^{(1)} + y^{(2)} \geq 10\}$.

Definition 1 (STL Syntax): Assume Π is the set of predicates and \mathcal{I} is any non-empty connected interval of $\mathbb{R}_{\geq 0}$. The set of all well-formed STL formulas is inductively defined as $\varphi ::= \top \mid \pi \mid \neg \phi \mid \phi_1 \vee \phi_2 \mid \bigcirc \phi \mid \phi_1 U_{\mathcal{I}} \phi_2$, where π is a predicate, \top is *true*, \bigcirc is Next, and $U_{\mathcal{I}}$ is the Until operator.

In this work, we will be using discrete time semantics of STL since we would like to be able to reason about the timing of samples and define events as falling or raising Boolean values using the next time operator (\bigcirc). For example, the

formula $\pi \wedge \bigcirc \neg \pi$ expresses an event (e.g., high to low). For STL formulas ψ, ϕ , we define $\psi \wedge \phi \equiv \neg(\neg\psi \vee \neg\phi)$, $\perp \equiv \neg\top$ (False), $\psi \rightarrow \phi \equiv \neg\psi \vee \phi$ (ψ Implies ϕ), $\diamond_I \psi \equiv \top U_I \psi$ (Eventually ψ), $\square_I \psi \equiv \neg \diamond_I \neg \psi$ (Always ψ), and $\psi R_I \phi \equiv \neg(\neg\psi U_I \neg\phi)$ (ψ Releases ϕ), using syntactic manipulation.

In our previous work [49], we proposed robust semantics for STL formulas. Robust semantics (or robustness metrics) provide a real-valued measure of satisfaction of a formula by a trace. In contrast, Boolean semantics just provide a *true* or *false* valuation. In more detail, given a trace σ of the system, its robustness w.r.t. a temporal property φ , denoted $\llbracket \varphi \rrbracket_d(\sigma)$, yields a positive value if σ satisfies φ and a negative value otherwise. Moreover, if the trace σ satisfies the specification ϕ , then the robust semantics evaluate to the radius of a neighborhood such that any other trace that remains within that neighborhood also satisfies the same specification. The same holds for traces that do not satisfy ϕ . In order to define neighborhoods for requirement satisfaction, we need to utilize metrics over the space of outputs, inputs and parameters.

Definition 2 (Metric): A metric on a set S is a positive function $d : S \times S \rightarrow \mathbb{R}_{\geq 0}$ such that

- I. $\forall s, s' \in S, d(s, s') = 0 \Leftrightarrow s = s'$
- II. $\forall s, s' \in S, d(s, s') = d(s', s)$
- III. $\forall s, s', s'' \in S, d(s, s'') \leq d(s, s') + d(s', s'')$

In this paper, all the experimental results were derived using the Euclidean metric, i.e., $d(s, s') = \|s - s'\|$; however, any other metric could be used. Using a metric, we can define a distance function that will capture how *robustly* a point belongs to a set. That is, the further away a point is from the boundary of the set, the more robust is its membership in that set.

Definition 3 (Signed Distance [50] §8): Let $s \in S$ be a point, $A \subseteq S$ be a set and d be a metric on S . Then, we define the Signed Distance from s to A to be

$$\text{Dist}_d(s, A) := \begin{cases} -\inf\{d(s, s') \mid s' \in A\} & \text{if } s \notin A \\ \inf\{d(s, s') \mid s' \notin A\} & \text{if } s \in A \end{cases}$$

That is, the signed distance is positive if the point is in the set and negative otherwise.

Definition 4 (STL Robust Semantics): Given a metric d , trace σ , and $\mathcal{O} : \Pi \rightarrow 2^{\mathcal{Y} \times \mathcal{U} \times \mathcal{P}}$, the robust semantics of any formula ϕ w.r.t σ at time instance $i \in \mathbb{N}$ is defined as:

$$\begin{aligned} \llbracket \top \rrbracket_d(\sigma, i) &:= +\infty \\ \llbracket \pi \rrbracket_d(\sigma, i) &:= \text{Dist}_d([y_i \ u_i \ p_i]^T, \mathcal{O}(\pi)) \\ \llbracket \neg \phi \rrbracket_d(\sigma, i) &:= -\llbracket \phi \rrbracket_d(\sigma, i) \\ \llbracket \phi_1 \vee \phi_2 \rrbracket_d(\sigma, i) &:= \max(\llbracket \phi_1 \rrbracket_d(\sigma, i), \llbracket \phi_2 \rrbracket_d(\sigma, i)) \\ \llbracket \bigcirc \phi \rrbracket_d(\sigma, i) &:= \begin{cases} \llbracket \phi \rrbracket_d(\sigma, i+1) & \text{if } i+1 \in \mathbb{N} \\ -\infty & \text{otherwise} \end{cases} \\ \llbracket \phi_1 U_I \phi_2 \rrbracket_d(\sigma, i) &:= \max_{j \text{ s.t. } (t_j - t_i) \in \mathcal{I}} \left(\min(\llbracket \phi_2 \rrbracket_d(\sigma, j), \right. \\ &\quad \left. \min_{i \leq k < j} \llbracket \phi_1 \rrbracket_d(\sigma, k)) \right) \end{aligned}$$

The value $\llbracket \phi \rrbracket_d(\sigma, 0)$ is referred to as the *robustness* with which σ satisfies ϕ . For convenience, we just write $\llbracket \phi \rrbracket_d(\sigma)$ when $i = 0$. As proved in [49], a trace σ satisfies an STL

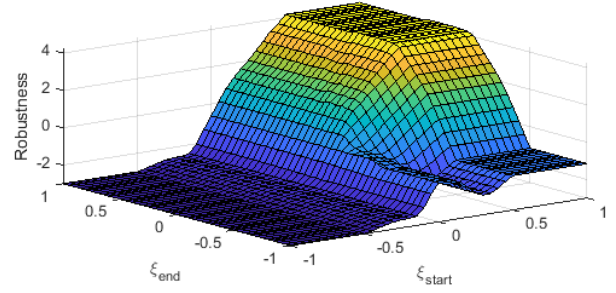


Fig. 4: The resulting robustness landscape (heatmap) for specification φ in Example 2. The initial positions for the two vehicles and the input μ are fixed. The input signals for ξ are generated by linear interpolation between two input values: ξ_{start} at time 0 and ξ_{end} at time 10.

formula ϕ (denoted by $\sigma \models \phi$), if $\llbracket \phi \rrbracket_d(\sigma) > 0$. On the other hand, a trace σ' does not satisfy ϕ (denoted by $\sigma' \not\models \phi$), if $\llbracket \phi \rrbracket_d(\sigma') < 0$. An overview of the algorithms that can be used to compute $\llbracket \varphi \rrbracket_d$ is provided in [48].

Example 2 (Continued from Example 1): Let's assume that we would like to check on the output traces of Fig. 3 the simple safety requirement: *after 5 time units, the distance between the adversarial and ego vehicles should always be greater than 0*. This requirement is captured by the STL specification $\varphi = \square_{[5, \infty)}(y^{(a)} - y^{(e)} > 0)$. Then, $\llbracket \varphi \rrbracket_d(\sigma) = -1.0841$, which means that the requirement is not satisfied. Moreover, the value -1.0841 corresponds to time 5.2446, which is the time of the worst violation of ϕ_1 .

B. Robustness-Guided Falsification

The robustness metric can be viewed as a fitness function that indicates the degree to which individual simulations of the system satisfy the requirement φ (positive values indicate that the simulation satisfies φ). Therefore, for a given system Σ and a given requirement φ , the verification problem is to ensure that for all $\sigma \in \mathcal{L}(\Sigma)$, $\llbracket \varphi \rrbracket_d(\sigma) > 0$.

Let φ be a given STL property that the system is expected to satisfy. The robustness metric $\llbracket \varphi \rrbracket_d$ maps each simulation trace σ to a real number r (see Fig. 4 for an example). Ideally, for the STL verification problem, we would like to prove that $\inf_{\sigma \in \mathcal{L}(\Sigma)} \llbracket \varphi \rrbracket_d(\sigma) > \varepsilon > 0$ where ε is a desired robustness threshold. Unfortunately, in general, the problem is not algorithmically solvable [51], that is, there does not exist an algorithm that can solve the problem. Hence, instead of trying to prove that the property holds on the system, we will try to demonstrate that it does not hold on the system when the system is unsafe. In other words, we are searching for a trajectory (trace) which *falsifies* the requirement (i.e., a trace that demonstrates that the specification is false). This is the topic of the next section.

C. Falsification and Critical System Behaviors

In this work, we focus on the task of identifying critical system behaviors, including falsifying traces. To identify falsifying system behaviors, we leverage existing work on *falsification*, which is the process of identifying system traces σ that

do not satisfy a given specification φ . The STL falsification problem is defined as: Find $\sigma \in \mathcal{L}(\Sigma)$ s.t. $\llbracket \varphi \rrbracket_d(\sigma) < 0$. One successful approach in addressing the falsification problem is to pose it as a global non-linear optimization problem:

$$\sigma^* = \arg \min_{\sigma \in \mathcal{L}(\Sigma)} \llbracket \varphi \rrbracket_d(\sigma). \quad (1)$$

If the global optimization algorithm converges to some local minimizer $\tilde{\sigma}$ such that $\llbracket \varphi \rrbracket_d(\tilde{\sigma}) < 0$, then a counterexample (adversarial sample) has been identified, which can be used for debugging (or for training). Considering the robustness heatmap in Fig. 4, any point with robustness below 0 would be a counterexample for the specification of Example 2. In order to solve this non-linear non-convex optimization problem, a number of stochastic search optimization methods can be applied (e.g., [22] – for an overview see [52], [53]). We leverage existing falsification methods to identify falsifying examples the autonomous driving system.

D. Covering Arrays

In software systems, there can often be a large number of discrete input parameters that affect the execution path of a program and its outputs. The possible combinations of input values can grow exponentially with the number of parameters. Hence, exhaustive testing on the input space becomes impractical for fairly large systems. A fault in such a system with k parameters may be caused by a specific combination of t parameters, where $1 \leq t \leq k$. One best-effort approach to testing is to make sure that all combinations of any t -sized subset (i.e., all t -way combinations) of the inputs are tested.

A *covering array* is a minimal number of test cases such that any t -way combination of test parameters exist in the list [21]. Covering arrays are generated using optimization-based algorithms with the goal of minimizing the number of test cases. We denote a t -way covering array on k parameters by $CA(t, k, (v_1, \dots, v_k))$, where v_i is the number of possible values for the i^{th} parameter. The size of the covering array increases with increasing t , and it becomes an exhaustive list of all combinations when $t = k$. Here, t is considered as the *strength* of the covering array. In practice, t can be chosen such that the generated tests fit into the testing budget. Empirical studies on real-world examples show that more than 90 percent of the software failures can be found by testing 2 to 4-way combinations of inputs [54].

Despite the t -way combinatorial coverage guaranteed by covering arrays, a fault in the system possibly may arise as a result of a combination of a number of parameters larger than t . Hence, covering arrays are typically used to supplement additional testing techniques, like uniform random sampling (fuzzing). We consider that because of the nature of the training data or the network structure, NN-based object detection algorithms may be sensitive to a certain combination of properties of the objects in the scene. Fig. 5 shows outputs of a DNN-based object detection and classification algorithm for 4 different combinations of vehicle type, vehicle color and pedestrian pants color while all other parameters like position and orientation of the objects are the same. In a comparison between configurations (a) and (b), the vehicle type does not

change but the vehicle and pedestrian pants colors change from blue to white. While both the car and the pedestrian are detected in configuration (a), the pedestrian is detected but the car is not detected in configuration (b); however, in a comparison between configurations (b) and (d), if we fix the vehicle and pedestrian pants colors to be white but change the vehicle type, then the car is detected but the pedestrian is not detected. We can also see that the size of the detection box is different between configurations (c) and (d), for which the vehicle type is the same but the vehicle and pedestrian pants colors are different. Our observation is that the characterization of the errors is generally not as simple as saying that all white colored cars are not detected. Instead, the errors arise from some combination of subsets of discrete parameters. Because of this combinatorial aspect of the problem, covering arrays is a good fit to test DNN-based object detection and classification algorithms. In Sec. V, we describe how Sim-ATAV combines covering arrays to explore discrete and discretized parameters with falsification on continuous parameters.

IV. REQUIREMENTS

In this section, we provide five STL requirements intended for the autonomous driving system. Each requirement is used to target specific aspects of safety and performance. Also, we describe how analysis results related to each of the requirements can be used to enhance either the controller design or a testing phase of the development process.

A. STL Requirements

This section describes each of the requirements that we use in the sequel to evaluate an ADS design with our virtual framework. We provide these requirements to illustrate how STL can be used to describe four different types of behavior expectations for an ADS: *system-level* safety, *subsystem-level*

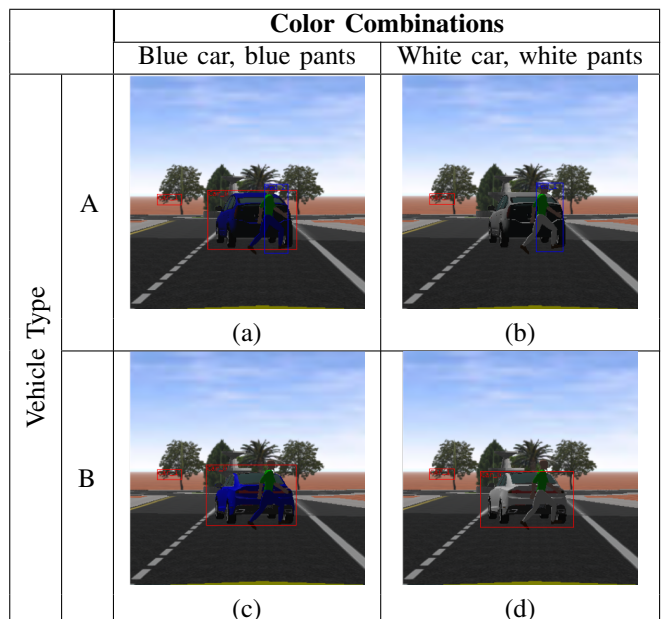


Fig. 5: Specific configurations impacting DNN performance.

performance, *subsystem-to-system* safety, and system-level performance (driving comfort) requirements. In the following, when an object i and a sensor s are clear from the context, we will drop the indices from the notation.

Requirement R1: The ego vehicle should not collide with an object.

This requirement is an example of a *system-level safety requirement*. It is used to ensure that the ego vehicle does not collide with any object in the environment. Behaviors that do not satisfy this requirement correspond to unsafe performance by the AV. These cases are valuable to identify in simulation, as they can be communicated back to the control designers so that the control algorithms can be improved.

The formal requirement in STL is:

$$R1_i = \Box(\neg\pi_{i,coll})$$

where $\pi_{i,coll} = dist(i, ego) < \epsilon_{dist}$.

In the above specification, i corresponds to an object in the environment, such as an agent vehicle or a pedestrian. $dist(i, ego)$ gives the minimum Euclidean distance between the boundaries of the Ego vehicle and the boundaries of object i . The specification basically indicates that the Ego vehicle should not collide with object i .

In practice, we consider a separate requirement for each object in the environment and all of them are checked conjunctively, i.e., if there are M_o objects, then $R1 = \bigwedge_{i=1}^{M_o} R1_i$.

Remark 2: As we have indicated in [30], [31], $R1$ can be too restrictive and pessimistic for an adversarial testing environment. Namely, if the requirement simply states “do not collide with a moving object”, then the test engine will attempt to generate agent trajectories that purposefully try to collide with the ego vehicle. If the adversaries are powerful enough, then the ego vehicle cannot avoid such collision cases. Hence, depending on the test scenario it may be necessary to enforce $R1$ only for static objects, or impose additional assumptions formalizing when the ego vehicle is supposed to be able avoid collisions, e.g., [55], [28].

Requirement R2: Sensor s should detect visible obstacles within t_1 time units.

This requirement is an example of a *subsystem-level requirement*. This particular example can be considered as a requirement on the sensor or perception subsystems. The requirement indicates that the perception system or a specific sensor s should not fail to detect an object for an excessive amount of time, i.e., more than t_1 time.

The requirement is as follows.

$$R2_{i,s} = \Box((W(i,s) \wedge \neg D(i,s)) \implies \Diamond_{[0,t_1]}(D(i,s) \vee \neg W(i,s)))$$

Here, $W(i,s)$ denotes that object i is physically visible to sensor s . For our framework, $s \in \{CCD, lidar, radar, combined\}$, where *combined* represents the total perception system, that is, fusion of all available sensors. The predicate $D(i,s)$ evaluates to true when sensor s detects object i . A (non-unique) description of this requirement in natural language is “it is always true that for any time when object i is visible and not detected by sensor

s , then there exists a time instant between 0 and t_1 that object i is either detected or it should be invisible to the sensor”.

Requirement R3: Localization errors should not be too large for too long.

This requirement is another sensor-level requirement and specifies that the localization of an object that is based on a particular sensor should provide sufficient accuracy, within an adequate time after the object becomes visible to the sensor:

$$R3_{i,s} = \Box((W(i,s) \wedge (\neg D(i,s) \vee E(i,s) > \epsilon_{err})) \implies \Diamond_{[0,t_1]}(\neg W(i,s) \vee (D(i,s) \wedge E(i,s) < \epsilon_{err})))$$

In $R3_{i,s}$, $E(i,s)$ is the difference between object i 's location and its location as estimated using information from sensor s . The constant ϵ_{err} is a threshold on the acceptable error between the actual position of i and its estimated position.

To understand the requirement, consider the situation where either an object is not detected (i.e., $\neg D(i,s)$) or there is a large error in the localization of the object (i.e., $E(i,s) > \epsilon_{err}$), we refer to this case as “poor detection” of the object. We can interpret the requirement as follows: “it is always true that whenever object i is visible to sensor s and is poorly detected by sensor s , then there exists an instant, within a time period from 0 to t_1 , that either object i is invisible to sensor s or the object is detected and the localization error is small, as computed using information from sensor s ”. This requirement basically limits the amount of time the sensor error can be greater than a given threshold.

Requirement R4: A sensor-related fault should not lead to a system-level fault.

This is an example of a *subsystem-to-system requirement*. This requirement relates sensor-level behaviors to system-level behaviors. The purpose is to isolate behaviors where a sensor fault results in a collision. The expectation is that the system as a whole should be robust to failure of a single sensor:

$$R4_{i,s} = \Box\neg\left(\Box_{[0,t_1]}(\neg\pi_{i,coll} \wedge W(i,s) \wedge (\neg D(i,s) \vee E(i,s) > \epsilon_{err})) \wedge \Diamond_{(t_1,t_2]}\pi_{i,coll}\right)$$

The above requirement designates that there should not be a period of t_1 where a visible object is not accurately detected and no collision occurs, followed immediately by a period of length $t_2 - t_1$ that contains a collision. In other words, the requirement indicates that a system level fault (collision) should not occur within a short time after a sensor fault. A behavior that violates this requirement does not necessarily indicate that the sensor fault *caused* the system fault, but it suggests a correlation, as it points to a behavior wherein the system fault occurs a short time after the sensor fault. Providing behavior examples that violate this requirement can help to pinpoint the cause of system-level faults.

Requirement R5: The vehicle should not do excessive braking unnecessarily or too often.

This is a *system-level performance (driving comfort)* requirement, in that it requires that the system not brake unnecessarily or too often, thereby causing discomfort for the passengers:

$$R5 = \square \left(\neg \square_{[0,t1]} (B \wedge \neg FC) \wedge \neg (B^\downarrow \wedge \diamond_{(0,t2)} (B^\downarrow \wedge \diamond_{(0,t2)} B^\downarrow)) \right),$$

Here, FC is a variable that is true when the Ego vehicle is estimated to collide in the future with another object in the environment, based on a simplified model of future behaviors. The simplified model that we use for future trajectory estimation is the Constant Turn Rate and Velocity (CTRV) model [56]. The proposition B represents that the amount of braking force applied by the controller exceeds half of the available braking force. Finally, $B^\downarrow = B \wedge \bigcirc \neg B$ represents the event of releasing the brake, i.e., a *true* value of B followed by a *false* value in the next sample.

To understand the meaning of requirement $R5$, consider the following part of the requirement:

$$\square \left(\neg \square_{[0,t1]} (B \wedge \neg FC) \right),$$

which requires that the system not apply excessive braking for more than a specific amount of time ($t1$) while there is no collision predicted. This essentially stipulates that the system should not unnecessarily brake for a prolonged amount of time. Next, consider the second part of requirement $R5$:

$$\square \left(\neg (B^\downarrow \wedge \diamond_{(0,t2)} (B^\downarrow \wedge \diamond_{(0,t2)} B^\downarrow)) \right),$$

which indicates that there should not be an “on-off” behavior, followed by another “on-off” behavior, followed by a third “on-off” behavior, with less than $t2$ between each other. This essentially requires that the brakes not be applied and released too often. Thus, this is a riding comfort requirement.

B. Development Process Support

We describe how requirements $R1$ through $R5$ can be used to support both the controller design and testing phases of the development process. For all of the requirements, any detected violation (falsification) should be linked back to the conditions that caused the violation.

Consider the first requirement, $R1$, “The vehicle should not collide with an object”: if the vehicle does collide with an object, then we would go back and see what conditions caused such an event, for example, whether the vehicle speed trace exhibited an anomaly or whether the vehicle was moving erratically. Testing for collision avoidance is well established in the field of ADAS. Often inflatable and other destructible targets are employed; for example, see [57] and Fig. 6.

In the requirement “Sensor should detect visible obstacles”, we focus on the detection of an obstacle as operational imperative. If the sensor fails to detect the object within a time interval, then the requirement is violated. This is essentially a sensor-level requirement (visible but not detected), and test engineers can set a real-world experiment to verify it relatively easy because it is decoupled from others (one-term inequality, sensor by sensor).

The requirement “Localization error should not be too high for too long” is important to verify (falsify) for both ego-location and position identification of other agents in the



Fig. 6: Robotic pedestrian surrogate target with a Toyota autonomous vehicle.

environment. Placing an ego-vehicle in the correct pose on the road is usually not achieved by simply relying on GPS signal processing, due to the GPS tendency to “jump” unpredictably, but instead by estimating and dynamically refining the pose through landmark observations, such as road edges, vertical elements such as light poles, and signs. Assuming that the ego-vehicle localization is done with sufficient accuracy, the remaining task of localization is to make sure that the location of other agents, especially those in the planned path of the ego vehicle, are estimated with sufficient accuracy. Often a grid-based representation centered on the ego-vehicle is employed (e.g., [58]). Estimating $E(i, s)$ in $R3$ is not trivial, but practical approaches exist that can be used by test engineers (e.g., [59]).

The requirement “A sensor-related fault should not lead to a system-level fault” is a form of robustness requirement. This is similar to a requirement that the system should have no “single point of failure”, which enforces that the failure of any single component will not cause the system to fail (for example, see [60]). We make an important clarification which is practical but limiting in scope: no failure should occur within the specified (short) time after the fault. Test engineers could readily use examples of behavior provided in the course of falsifying this requirement.

Lastly, the requirement “The vehicle should not brake too often” is an example of a possible set of requirements designed to establish how comfortable the ride in the vehicle is. It is known that autonomous vehicles could induce motion sickness in passengers if the vehicle control system does not comply with human physiology [61], [62]. A better requirement may well be developed using fuzzy set theory and further refined for a specific target group of passengers (e.g., elderly people). An alternative requirement could be defined by counting the number of occurrences of an event within a total time period, instead of relating one occurrence to another. Such a requirement can be defined as a Timed Propositional Temporal Logic (TPTL) specification. TPTL is a generalization of STL which is also supported in our framework [63].

V. SIM-ATAV FRAMEWORK

We describe Sim-ATAV, a framework for performing testing and analysis of autonomous driving systems in a virtual environment. The framework is publicly available as an add-on to S-TaLiRo [24]:

<https://sites.google.com/a/asu.edu/s-taliro/s-taliro/sim-atav>

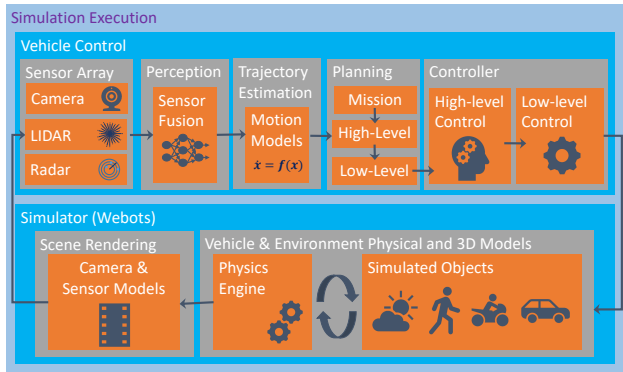


Fig. 7: Overview of the simulation environment.

The simulation environment used in Sim-ATAV is based on the open source simulator Webots [64] and includes a vehicle perception system, a vehicle controller, and a model of the physical environment. The perception system processes data from three sensor systems: CCD camera images, lidar, and radar. The framework uses freely available and low cost tools and can be run on a standard desktop PC. Later, we demonstrate how Sim-ATAV can be used to implement traditional testing approaches, as well as advanced automated testing approaches that were not previously possible using existing frameworks.

Fig. 7 shows an overview of the simulation environment. The environment consists of a Simulator and a Vehicle Control system. The Simulator contains models of the ego vehicle, agents, and other objects in the environment (e.g., roads, buildings). The Simulator outputs sensor data to the Vehicle Control system. The sensor data includes representations of CCD camera, lidar, and radar data. Simple models of the sensors are used to produce the sensor data. For example, synthetic CCD camera images are rendered by the Simulation system, as if they came from a camera mounted on the front of the ego vehicle. The Vehicle Control system contains models of the Perception System, which performs sensor data processing and sensor fusion. The Controller uses the output of the Perception System to make decisions about how to actuate the AV system. Actuation commands are sent from the Controller to the Simulator.

Simulations proceed iteratively. At each instant, sensor data is processed by the Vehicle Control, which then makes an actuation decision. The actuation decision is then transmitted back to the Simulator, which uses the actuation commands to update the physics for the next time instant. This process is repeated until a designated time limit has been reached.

The Vehicle Control system is implemented in Python. We use simplified algorithms to implement the subsystems of the vehicle control, which is sufficient in this case, as the purpose of this investigation is to evaluate new *testing* methodologies and not to evaluate a real AV control design; however, we note that it is straightforward to replace our algorithms with production versions to test real control designs.

To process CCD image data, we use a lightweight DNN, SqueezeDet, which performs object detection and classification [12]. SqueezeDet is implemented in TensorFlowTM[65],

and it outputs a list of object detection boxes with corresponding class probabilities. This network was originally trained on real image data from the KITTI dataset [11] to achieve accuracy comparable to the popular AlexNet classifier [17]. We further train this network on the virtual images generated in our framework. Fig. 8 shows an example output from SqueezeDet, based on a synthetic image produced by our simulator. The image shows two vehicles correctly detected and classified, along with a portion of a shadow that is incorrectly classified as a vehicle.



Fig. 8: Outputs from the SqueezeDet DNN, based on a synthesized camera image.

To process lidar point cloud data, we first cluster the received points based on their positions using the DBSCAN algorithm [66], [67]. Then, we estimate the existence and types of the objects based on how well the dimensions of the clusters match with the dimensions of expected object types such as pedestrians or cars. For estimating the object type in the received radar targets, we use the radar signal power. We implement a simple sensor fusion algorithm that relates and merges the object detections from camera, lidar, and radar with a simple logic. Our sensor fusion algorithm makes some rule-based decisions, such as if the object type can be recognized by the camera, discard the type estimations done by the lidar and radar, and on the other hand if radar or lidar is able to detect to position of the object, discard the position estimations computed by the camera. It also utilizes the expected current positions of previously detected objects. A simple implementation of an unscented Kalman filter is used to estimate the current and future trajectories of the objects using the CTRV (Constant Turn Rate and Velocity) model [68], [56].

Fig. 9 illustrates outputs from the sensor fusion system. In the figure, the solid yellow box in the middle represents the Ego vehicle. Yellow circles in front of the ego vehicle represent the estimated future trajectory of the Ego vehicle. Small white dots represent lidar point cloud data. The colored dots and rectangles represent detected objects, with their estimated orientation indicated with a white line in front of

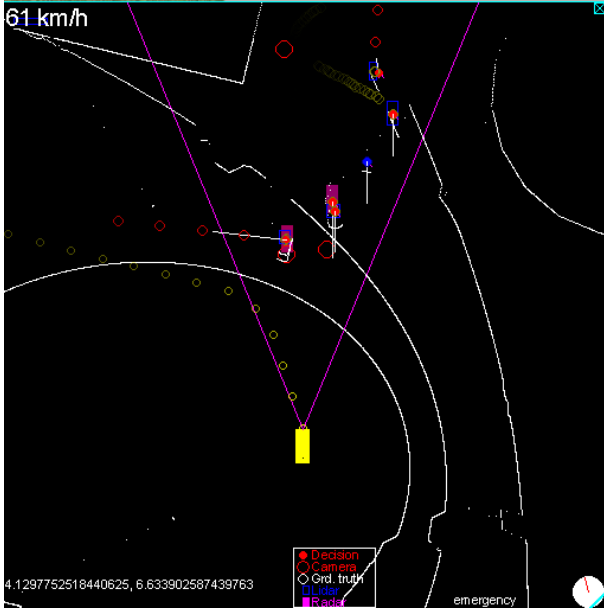


Fig. 9: Sensor fusion outputs.

them. Expected future positions of agent vehicles with respect to the ego vehicle are represented by red circles.

Our simple planner receives as inputs the high level target path and target speed, and the outputs of the sensor fusion and trajectory estimation modules. It assigns collision risk level to the target objects with a simple logic and outputs the risk assessments and a target speed, which depends on the target speed of the mission or other factors, such as the distance to a sharp turn ahead.

Our control algorithm implements simple path and speed tracking and *collision avoidance* features. The controller receives the outputs of the planner. When there is no collision risk, the controller drives the car with the target speed and on the target path. When a future collision with an object is predicted, it applies the brakes at a level proportional with the risk assigned to the object.

The environment modeling framework is implemented in Webots [64], a robotic simulation framework that models the physical behavior of robotic components, such as manipulators and wheeled robots, and can be configured to model autonomous driving scenarios. In addition to modeling the physics, a graphics engine is used to produce images of the scenarios. In Sim-ATAV, the images rendered by Webots are configured to correspond to the image data captured from a virtual camera that is attached to the front of a vehicle.

The process used by Sim-ATAV for test generation and execution for discrete and discretized continuous parameters is illustrated by the flowchart shown in Fig. 10-(a). Sim-ATAV first generates test cases that correspond to scenarios defined in the simulation environment using covering arrays as a combinatorial test generation approach. The scenario setup is communicated to the simulation interface using TCP/IP sockets. After a simulation is executed, the corresponding simulation trace is received via socket communication and evaluated using a cost function. Among all discrete test cases,

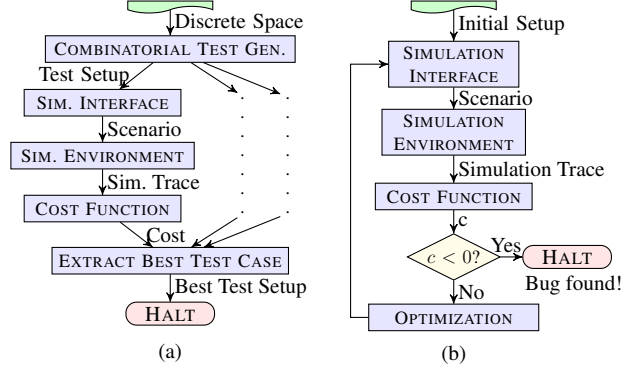


Fig. 10: Flowcharts illustrating the combinatorial testing (a) and falsification (b) approaches.

the most promising one is used as the initial test case for the falsification process shown in Fig. 10-(b). For falsification, the result obtained from the cost function is used in an optimization setting to generate the next scenario to be simulated. For this purpose, we used S-TaLiRo [24], which is a MATLAB® toolbox for falsification of CPSs. Similar tools, such as Breach [69], can also be used in our framework for the same purpose.

VI. TESTING APPLICATION

In this section, we present an evaluation of our Sim-ATAV framework using three separate driving scenarios. The scenarios are selected to be both challenging for the autonomous driving system and also analogous to plausible driving scenarios experienced in real-world situations. In general, the system designers will need to identify crucial driving scenarios, based on intuition about challenging situations, from the perspective of the autonomous vehicle control system. A thorough simulation-based testing approach will include a wide array of scenarios that exemplify critical driving situations.

For each of the following scenarios, we consider a subset of the requirements presented in Sec. IV and describe how to use the results to enhance the development process. We conclude the section with a summary of the results.

Scenario 1

The first scene that we consider is a straightaway section of a two-lane road, as illustrated in Fig. 11. Several cars are parked on the right-hand side of the road, and a pedestrian is jay-walking in front of one of the cars, passing in front of the Ego car from right to left. We call this driving scenario model M_1 . The scenario simulates a similar setup to the Euro NCAP Vulnerable Road User (VRU) protection test protocols [4].

Several aspects of the driving scenario are parameterized, meaning that their values are fixed for any given simulation by appropriately selecting the model parameters. The parameters and initial conditions that we use for this scenario are:

- Initial speed of the Ego vehicle: $[10, 30]m/s$;
- Lateral position of Ego w.r.t its lane center: $[-0.8, 0.8]m$;
- Walking speed of the pedestrian: $[1.5, 6]m/s$;
- The model of Agent car, which is next to the pedestrian: *from 5 different vehicle models*;

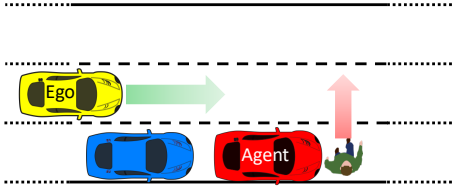


Fig. 11: Overview of the scenario 1.

- R, G, B values for the colors of Agent car: $[0, 1]$;
- R, G, B values for the pedestrian's shirt and pants: $[0, 1]$.

We chose the parameterized aspects of the scenarios such that their specific combinations would be challenging to a DNN-based pedestrian detection system that relies on CCD camera images, *e.g.*, some combinations of agent vehicle shape and color with pedestrian clothing colors may be challenging to an object detection DNN. We also chose some of the parameter ranges, *i.e.*, the vehicle and pedestrian speeds, so that the scenario is physically challenging for the brake performance. The Ego vehicle is longitudinally placed such that the time-to-collision with the pedestrian will be $1.4s$ when the pedestrian is exactly in front of the Ego vehicle (influenced by the Euro NCAP, VRU protection test protocols [4]).

We evaluate Model M_1 against three of the requirements from Sec. IV: $R1$, $R2$, and $R4$. These include the system-level requirement, the sensor-level requirements, and the sensor-to-system-level requirement. We use this collection of requirements for Model M_1 to demonstrate how we can automatically identify each type of behavior using our framework.

Scenario 2

The next scenario involves a left turn maneuver by the ego vehicle in a controlled intersection, as illustrated in Fig. 12. An agent vehicle (*Agent 1*) in the opposing lane unexpectedly passes through the intersection, against a red light, potentially causing a collision with the Ego vehicle. There is also another agent car (*Agent 2*), which is making a legal left turn from the opposing lane. It is incumbent on the Ego vehicle to take action to avoid colliding with the agent vehicles. We call the model of this scenario M_2 .

For this experiment, we choose parameters such that the position of Agent 2, or trajectory followed by Agent 1, in combination with the behavior of the Ego, may result in poor performance from the sensor processing or trajectory estimation systems. For this scenario the search space is:

- Ego vehicle initial speed: $[20, 30]m/s$;
- Ego vehicle initial distance to the intersection: $[80, 160]m$;
- Agent 1 initial distance to the intersection: $[50, 100]m$;
- Agent 1 target speed (initial, when approaching the intersection, when inside the intersection): $[10, 35]m/s$;
- Agent 1 lateral position w.r.t its lane center (initial, when approaching the intersection, inside the intersection): $[-0.75, 0.75]m$;
- Agent 2 lateral position w.r.t its lane center: $[-0.75, 0.75]m$;
- Agent 2 speed: $[3, 15]m/s$;

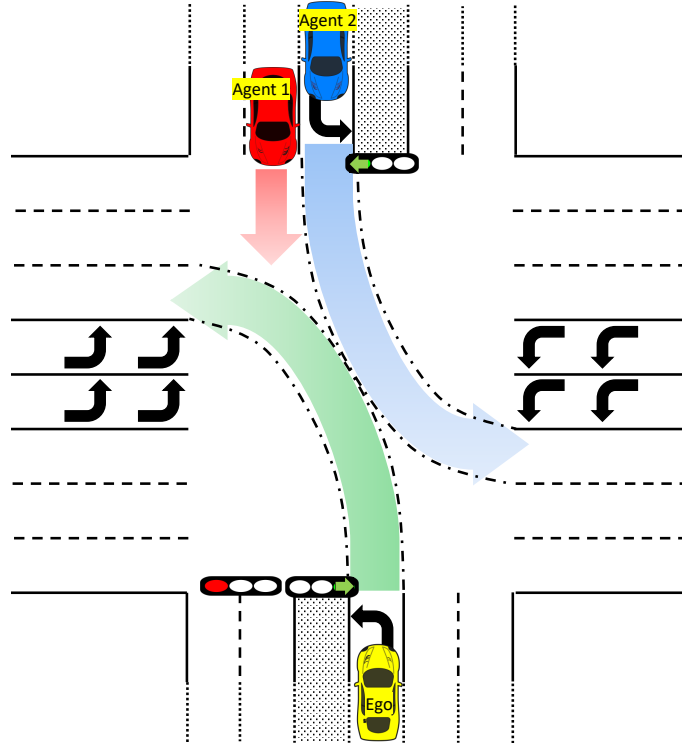


Fig. 12: Overview of the scenario 2.

- Agent 2 initial distance to the intersection: $[50, 100]m$.

We evaluate Model M_2 against requirement $R4$. The idea in using the sensor-to-system-level requirement is that it is relatively easy, in general, to find behaviors that result in a collision for Model M_2 , but many collision cases are not interesting for the designers. This could be because, for example, the agent car is moving too quickly for the ego vehicle to avoid. This would be a behavior that is not necessarily caused by any specific incorrect behavior on the part of the ego vehicle. Instead, we use $R4$ to identify behaviors where there is a collision that is directly correlated to unacceptable performance from the sensor processing system. These are cases where the sensor data processing or future trajectory estimation system is at fault for the collision. Such cases can be easily used to debug specific aspects of the ego vehicle control algorithms.

Scenario 3

In this last scenario, the ego vehicle is making a left turn through an intersection, while an agent vehicle in the opposing lane is also making a left turn. This scene is similar to the Scenario 2, as depicted in Fig. 12, except that Agent 1 is not present in this scenario, only Agent 2, which we refer to as the agent vehicle for this scenario. If both ego and agent vehicles are not accurately regulating their trajectories during this maneuver, a collision may occur. We call the model of this scenario M_3 .

In this scenario, we search over target trajectories of the ego and agent vehicles. Below are the parameters that we use:

- Ego vehicle initial speed: $[15, 30]m/s$;

- Ego target lateral position w.r.t its lane center when entering the intersection: $[-1, 1]m$, longitudinal starting position of the left turn w.r.t. intersection entry point: $[-1, 10]m$, target lateral position w.r.t its lane center when exiting the the intersection: $[-1, 1]m$, end position of the left turn w.r.t. intersection exit point: $[-1, 10]m$;
- Agent vehicle speed: $[15, 30]m/s$, target lateral position w.r.t its lane center when entering the intersection: $[-1, 1]m$, longitudinal starting position of the left turn w.r.t. intersection entry point: $[-1, 10]m$, target lateral position w.r.t its lane center when exiting the the intersection: $[-1, 1]m$, and the end position of the left turn w.r.t. intersection exit point: $[-1, 10]m$.

We evaluate Model M_3 against requirement $R5$. The purpose of considering the performance requirement $R5$ in this case, is that scenario M_3 is difficult to falsify. That is, due to the specific parameter ranges selected for the scenario, it is unlikely that the ego vehicle will collide with the agent vehicle. Instead, in this case, we are interested in identifying situations where the emergency braking system unnecessarily decelerates the ego vehicle, causing unacceptable performance, from a ride-quality perspective. The scenario can easily lead to unnecessary braking, as the ego and agent vehicles momentarily move toward each other during their left turn maneuvers, which can cause the emergency braking algorithm to decide, incorrectly, that a collision is imminent. Such cases can be a useful feedback to designers, as they can highlight behaviors that are too conservative at the expense of ride quality.

Summary of Test Results

We present results from experiments demonstrating the application of our framework to the scenarios and requirements described above. Table I summarizes the results. For each case study, Table I indicates the requirements used to test each model, the testing approach used, the set of active sensors used, and a summary of the results. We discuss the test generation outcomes in detail below.

Covering array and falsification on Model M_1 : In our previous work [23], we proposed and studied the effectiveness of a testing approach that first uses covering arrays to discover critical regions, based on a set of discrete parameters, then uses those results as the initial points for robustness guided falsification. Here, we apply that approach on model M_1 for 3 different requirements, $R1$, $R2$ and $R4$. In model M_1 , we focus on the camera sensor and DNN-based object detection and classification algorithm. Because of this, most of our parameters are colors of pedestrian clothing and the agent vehicle, to which an object detection system may be sensitive to as described in Sec. IV.

We create a mixed-strength covering array from a discretization of the parameterized variables. We choose a 3-way covering of the parameters *ego vehicle initial speed*, *ego vehicle lateral position*, *pedestrian speed*, *agent vehicle model* as they are intuitively the most critical parameters for the collision avoidance performance, and a 2-way covering of the other parameters. These settings result in 195 covering array tests generated by the ACTS tool [54]. Depending on the

time budget for the testing, the strength of the covering array parameters can be changed, which can drastically increase or decrease the number of discrete test cases in the generated covering array. We first execute the resulting 195 covering array tests and collect simulation trajectories. Then, we compute the *robustness* values for those trajectories, with respect to the requirements $R1$, $R2$ and $R4$. Finally, for each requirement, starting from the case with the smallest positive robustness value, we try to find as many additional falsifications as possible, within a maximum of 300 extra simulations, by using a falsification approach that uses simulated annealing to perform the optimization.

For requirement $R1$, 67 cases were falsified from the covering array tests (i.e., 67 of the 195 cases did not satisfy $R1$). Starting from 7 of the remaining (non-falsifying) cases from the covering array tests, 5 additional falsifying cases were discovered using falsification. For requirement $R2$, 65 cases were falsified from the covering array cases, with an additional 8 cases discovered during the falsification step. For requirement $R4$, 67 cases were falsified during the covering array step, with 12 more cases discovered during the falsification step.

These results demonstrate that we can automatically identify test cases that violate specific sensor-level, system-level, and sensor-to-system level requirements. These test cases can be fed back to the designers to improve the perception or control design or can be used as guidance to identify challenging scenarios to be used during the testing phase.

Analysis of robustness values on the falsification of Model M_2 : The *robustness* value, which is described in Sec. III, for a trajectory with respect to the requirement is automatically computed in Sim-ATAV. This computation is performed by the S-TaLiRo tool [24] and is used to guide the test cases towards a falsification.

We use the results of falsification on Model M_2 to show, in Fig. 13, how the robustness value changes over time and finally becomes negative, which indicates falsification of the requirement. In this case, Sim-ATAV was able to find a falsifying example in 58 simulations. Because the cost function gradients are not computable, we use a stochastic global optimization technique, Simulated Annealing (SA). The blue line shows the robustness value for each simulation. We can observe that the robustness value per simulation run is not monotonic. This is due to the stochastic nature of the optimizer; however, the achieved minimum robustness up to the current simulation is a non-increasing function, which shows the best robustness achieved after each simulation. As soon as the framework finds a test case that causes a negative robustness value, it stops the search and reports the falsifying example.

Fig. 14 shows images from the simulation execution of a falsifying example for model M_2 with respect to the requirement $R2$. Between the time corresponding to Fig. 14-(a) to Fig. 14-(b), the red car approaching from the opposite side is driving on a path such that there will be a future collision with the Ego vehicle. However, due to incorrect localization of the agent vehicle, the Ego vehicle is not able to correctly predict the future trajectory of the agent vehicle, and so it does not predict a collision. Hence, it continues without taking action

	Model M_1			Model M_2	Model M_3
Requirement	$R1$	$R2$	$R4$	$R4$	$R5$
Testing Modality	CA+ Falsification	CA+ Falsification	CA+ Falsification	Falsification	Falsification
Active Sensors	CCD	CCD	CCD	CCD, Radar, LIDAR	CCD, LIDAR
Computation Time	CA: 2h, 10min.			2h, 3min.	9h, 40min.
	Fals.:3h, 33min.	3h, 35min.	3h 34min.		
No. Simulations	CA: 195			58	232
	Fals.:300	300	300		
Falsification Obtained	67 by CA + 5 by falsification	65 by CA + 8 by falsification	67 by CA + 12 by falsification	✓	✓
Application of Results	Lowest robustness cases used to create critical tests.			Falsifying cases relate to processing of specific sensor; aids in controller design improvement.	Poor performance cases used to improve controller design in modeling phase.

TABLE I: Results from autonomous driving tests using virtual framework.

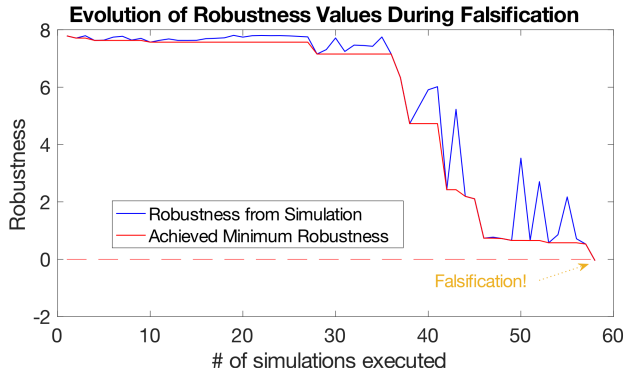


Fig. 13: Robustness guided falsification utilizes global optimization techniques to guide the test cases toward falsification.

to avoid the collision. Starting from the moment shown on Fig. 14-(c), the Ego vehicle predicts the collision and starts applying emergency braking; however, because it takes action too late, the Ego vehicle cannot avoid a collision with the agent vehicle, as shown in Fig. 14-(d).

We note that, even for cases that are non-falsifying, the robustness values are useful for the system designer’s analysis, as behaviors with low robustness value are “close to” violating the requirement and, therefore, correspond to cases that may require closer attention.

A visual analysis of a falsifying simulation trajectory from Model M_3 : As presented in Table I, Sim-ATAV was able to find a falsifying example for model M_3 with respect to the STL requirement $R5$ in 232 simulations. We present a visual analysis of the falsifying test result. Note that this analysis is done automatically in the framework, and corresponding satisfaction/falsification of the requirement is returned to the user, along with the *robustness* value that shows the signed distance to the boundary of satisfaction or falsification. The type of visual analysis we present here may be useful for the system designers to understand the reason behind the falsification (or satisfaction) of a requirement, which can be helpful for debugging or improving the design. For this analysis, we use the definitions and notation introduced in

Section IV.

Fig. 15 shows a part of the simulation trajectory of Model M_3 for a time window around the falsification instance, together with the corresponding logic evaluations of the predicates related to the subformulas in Requirement $R5$. In the top plot in Fig. 15, the red solid line is the estimated future minimum distance between the ego vehicle and Agent vehicle 1, with respect to the simulation time. This estimation is based on the ground truth information collected from the simulation and utilizes the CTRV model at each time step of the simulation to compute the collision estimate that is described in $R5$. For this example, we define the variable FC that is used in $R5$ as $(d_{f,min} < 0.5)$, where $d_{f,min}$ represents the expected minimum future distance. The dashed horizontal red line in Fig. 15 located at $0.5m$ is the threshold minimum future distance for a collision prediction. The values of $t1$ and $t2$ are respectively defined as 0.6 and 0.5 in this example. Since $d_{f,min}$ is never less than 0.5 in this case, the collision estimation variable FC , which is represented by the black solid line in the top plot, is always false.

The middle plot presents a similar evaluation for computing the variable B used in $R5$, which represents excessive braking. This evaluation uses the collected actual normalized brake power data, say br , from the simulation and computes the logical variable $B = (br > 0.5)$. The solid and dashed red lines represent br and the threshold value 0.5, respectively. The solid black line shows the value of B with respect to time. The bottom plot in Fig. 15 shows the value of the variable B^\downarrow that is defined for the requirement $R5$ with respect to the simulation time.

The first part of the requirement $R5$, which was defined as $\square(\neg\square_{[0,t1]}(B \wedge \neg FC))$ in Section IV, would evaluate to false if and only if there exists a time window of $t1$ such that B is always true and FC is always false. Focusing on the values of FC and B from the top two plots in Fig. 15, we can see that although FC is always false, because there is no time window of $t1 = 0.6s$ in which B is always true, the first part of the requirement evaluates to true. This means this execution of model M_3 satisfies the first part of the requirement $R5$.

The second part of the requirement $R5$, which is defined as $\square(\neg(B^\downarrow \wedge \diamond_{(0,t2]}(B^\downarrow \wedge \diamond_{(0,t2]}B^\downarrow)))$ evaluates to false if and only if there exists a series of three brake releases (B^\downarrow), such that one occurrence of B^\downarrow follows another within a $t2$ time window. As we see in the bottom plot of Fig. 15, at time 5.6s it is true that there exists B^\downarrow and it is also true that there exists another B^\downarrow within the time window of 0 to 0.5s following this moment (occurring at 5.85s). Hence, the inner $(B^\downarrow \wedge \diamond_{(0,t2]}B^\downarrow)$ inside the above formula evaluates to true at time 5.6s. If we call this event $e1$, the overall formula will evaluate to false if there exists an B^\downarrow that is followed by event $e1$ in a time window between 0 and $t2 = 0.5s$. This happens at time 5.46s, which is the moment that there exists an B^\downarrow followed by event $e1$ at $0.14 \in (0, 0.5]$, where the event $e1$ is defined as an B^\downarrow followed by another B^\downarrow within $t \in (0, 0.5]$. Hence, the second part of the requirement $R5$ evaluates to false, and as a result, $R5$ evaluates to false at time 5.46s, since it is a conjunction of parts 1 and 2. In other words, the system falsifies (does not satisfy) the requirement $R5$.

VII. CONCLUSIONS

We demonstrated a simulation-based adversarial test generation framework for autonomous vehicles. The framework works in a closed-loop fashion, where the system evolves in time with the feedback cycles from the autonomous vehicle’s controller. The framework includes models of lidar and radar sensor behaviors, as well as a model of the CCD camera sensor inputs. CCD camera images are rendered synthetically by our framework and processed using a pre-trained deep neural network (DNN). Using our framework, we demonstrated a new effective way of finding a critical vehicle behaviors by using 1) covering arrays to test combinations of discrete parameters and 2) simulated annealing to find corner-cases.

Future work will include using identified counterexamples to retrain and improve the DNN-based perception system, e.g., [70]. Additionally, the scene rendering will be made

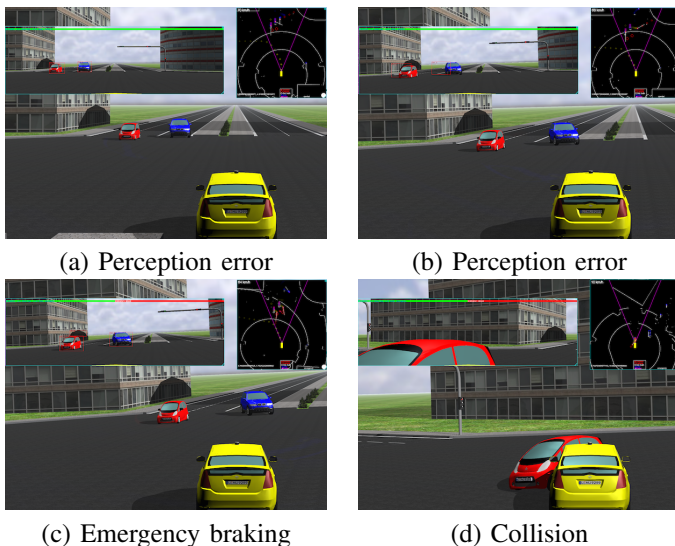


Fig. 14: Time-ordered images from the falsifying example on model M_1 .

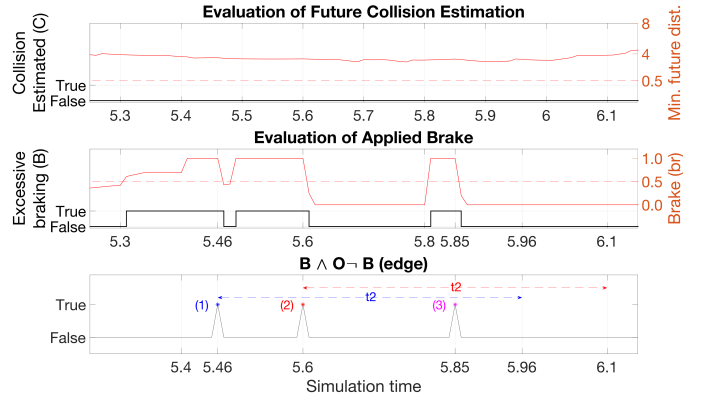


Fig. 15: Analysis of falsification for Model M_3 .

more realistic by using other scene rendering tools, such as those based on state-of-the-art game engines, e.g., [71]. Also, we note that the formal requirements that we considered were provided as an example of the type used when employing a requirements-driven development approach based on a temporal logic language, which is a formalism that may be unfamiliar to many test engineers. Future research will investigate ways to elicit formal requirements based for ADS using visual specification languages, e.g., [72].

REFERENCES

- [1] “NHTSA Federal Automated Vehicles Policy,” <https://www.transportation.gov/AV/federal-automated-vehicles-policy-september-2016>, accessed: 2018-09-11.
- [2] SAE, “Guidelines for safe on-road testing of sae level 3, 4, and 5 prototype automated driving systems (ADS), document j3018_201503,” https://www.sae.org/standards/content/j3018_201503/, accessed: 2018-09-12.
- [3] A. Christensen, A. Cunningham, J. Engelman, C. Green, C. Kawashima, S. Kiger, D. Prokhorov, L. Tellis, B. Wendling, and F. Barickman, “Key considerations in the development of driving automation systems,” in *24th Enhanced Safety of Vehicles Conferences*, 2015.
- [4] “Euro NCAP,” <https://www.euroncap.com/en>, accessed: 2018-09-11.
- [5] “Pegasus,” <https://www.pegasusprojekt.de/en/home>, accessed: 2018-09-11.
- [6] J. E. Stellet, M. R. Zofka, J. Schumacher, T. Schamm, F. Niewels, and J. M. Zöllner, “Testing of advanced driver assistance towards automated driving: A survey and taxonomy on existing approaches and open questions,” *2015 IEEE 18th International Conference on Intelligent Transportation Systems*, pp. 1455–1462, 2015.
- [7] M. R. Zofka, M. Essinger, T. Fleck, R. Kohlhaas, and J. M. Zöllner, “The sleepwalker framework: Verification and validation of autonomous vehicles by mixed reality lidar stimulation,” in *SIMPAR*. IEEE, 2018, pp. 151–157.
- [8] P. Koopman and M. Wagner, “Challenges in autonomous vehicle testing and validation,” *SAE Int. J. Trans. Safety*, vol. 4, pp. 15–24, 04 2016. [Online]. Available: <https://doi.org/10.4271/2016-01-0128>
- [9] W. Burgard, U. Franke, M.ENZweiler, and M. Trivedi, “The Mobile Revolution - Machine Intelligence for Autonomous Vehicles (Dagstuhl Seminar 15462),” *Dagstuhl Reports*, vol. 5, no. 11, pp. 62–70, 2016. [Online]. Available: <http://drops.dagstuhl.de/opus/volltexte/2016/5764>
- [10] W. Wachenfeld, P. Junietz, R. Wenzel, and H. Winner, “The worst-time-to-collision metric for situation identification,” in *2016 IEEE Intelligent Vehicles Symposium, IV 2016, Gotenburg, Sweden, June 19-22, 2016*, 2016, pp. 729–734.
- [11] A. Geiger, P. Lenz, and R. Urtasun, “Are we ready for autonomous driving? the kitti vision benchmark suite,” in *Computer Vision and Pattern Recognition (CVPR), 2012 IEEE Conference on*. IEEE, 2012.

- [12] B. Wu, F. Iandola, P. H. Jin, and K. Keutzer, "SqueezeDet: Unified, small, low power fully convolutional neural networks for real-time object detection for autonomous driving," in *IEEE Conference on Computer Vision and Pattern Recognition Workshops (CVPRW)*, 2017, pp. 446–454.
- [13] D. A. Pomerleau, "Alvin: An autonomous land vehicle in a neural network," in *Advances in neural information processing systems*, 1989, pp. 305–313.
- [14] C. Chen, A. Seff, A. Kornhauer, and J. Xiao, "DeepDriving: Learning affordance for direct perception in autonomous driving," in *Proceedings of the IEEE International Conference on Computer Vision*, 2015, pp. 2722–2730.
- [15] L. Chi and Y. Mu, "Deep steering: Learning end-to-end driving model from spatial and temporal visual cues," *arXiv preprint arXiv:1708.03798*, 2017.
- [16] M. Strickland, G. Fainekos, and H. B. Amor, "Deep predictive models for collision risk assessment in autonomous driving," in *IEEE International Conference on Robotics and Automation (ICRA)*, 2018.
- [17] A. Krizhevsky, I. Sutskever, and G. E. Hinton, "Imagenet classification with deep convolutional neural networks," in *Advances in neural information processing systems*, 2012, pp. 1097–1105.
- [18] Y. Tian, K. Pei, S. Jana, and B. Ray, "Deepest: Automated testing of deep-neural-network-driven autonomous cars," in *40th International Conference on Software Engineering (ICSE)*, 2018.
- [19] N. Papernot, P. McDaniel, I. Goodfellow, S. Jha, Z. B. Celik, and A. Swami, "Practical black-box attacks against machine learning," in *ACM Asia Conference on Computer and Communications Security*. ACM, 2017.
- [20] M. Wicker, X. Huang, and M. Kwiatkowska, "Feature-guided black-box safety testing of deep neural networks," in *International Conference on Tools and Algorithms for the Construction and Analysis of Systems (TACAS)*, ser. LNCS, vol. 10805, 2018, pp. 408–426.
- [21] A. Hartman, "Software and hardware testing using combinatorial covering suites," *Graph theory, combinatorics and algorithms*, vol. 34, pp. 237–266, 2005.
- [22] H. Abbas, G. E. Fainekos, S. Sankaranarayanan, F. Ivancic, and A. Gupta, "Probabilistic temporal logic falsification of cyber-physical systems," *ACM Transactions on Embedded Computing Systems*, vol. 12, no. s2, May 2013.
- [23] C. E. Tuncali, G. Fainekos, H. Ito, and J. Kapinski, "Simulation-based Adversarial Test Generation for Autonomous Vehicles with Machine Learning Components," in *IEEE Intelligent Vehicles Symposium (IV)*, 2018.
- [24] Y. S. R. Annapureddy, C. Liu, G. E. Fainekos, and S. Sankaranarayanan, "S-TaLiRo: A tool for temporal logic falsification for hybrid systems," in *Tools and algorithms for the construction and analysis of systems*, ser. LNCS, vol. 6605. Springer, 2011, pp. 254–257.
- [25] W. G. Najm, R. Ranganathan, G. Srinivasan, J. S. Toma, E. Swanson, and A. B. D. Smith, "Description of light-vehicle pre-crash scenarios for safety applications based on vehicle-to-vehicle communications," DOT HS 811 731, Tech. Rep., 2013.
- [26] D. Zhao, Y. Guo, and Y. J. Jia, "TrafficNet: An open naturalistic driving scenario library," in *20th IEEE International Conference on Intelligent Transportation Systems, ITSC*, 2017.
- [27] D. Zhao, H. Lam, H. Peng, S. Bao, D. J. LeBlanc, K. Nobukawa, and C. S. Pan, "Accelerated evaluation of automated vehicles safety in lane-change scenarios based on importance sampling techniques," *IEEE Transactions on Intelligent Transportation Systems*, vol. 18, no. 3, pp. 595–607, 2017.
- [28] S. M. Loos, A. Platzer, and L. Nistor, "Adaptive cruise control: Hybrid, distributed, and now formally verified," in *Formal Methods*, ser. LNCS, vol. 6664. Springer, 2011, pp. 42–56.
- [29] M. Althoff, D. Althoff, D. Wollherr, and M. Buss, "Safety verification of autonomous vehicles for coordinated evasive maneuvers," in *IEEE Intelligent Vehicles Symposium*, 2010.
- [30] C. E. Tuncali, T. P. Pavlic, and G. Fainekos, "Utilizing S-TaLiRo as an automatic test generation framework for autonomous vehicles," in *IEEE Intelligent Transportation Systems Conference*, 2016.
- [31] C. E. Tuncali and G. Fainekos, "Rapidly-exploring random trees for testing automated vehicles," in *IEEE Intelligent Transportation Systems Conference (ITSC)*, 2019.
- [32] M. Althoff and S. Lutz, "Automatic generation of safety-critical test scenarios for collision avoidance of road vehicles," in *IEEE Intelligent Vehicles Symposium (IV)*, 2018.
- [33] C. E. Tuncali, S. Yaghoubi, T. P. Pavlic, and G. Fainekos, "Functional gradient descent optimization for automatic test case generation for vehicle controllers," in *IEEE International Conference on Automation Science and Engineering*, 2017.
- [34] M. O'Kelly, H. Abbas, and R. Mangharam, "Computer-aided design for safe autonomous vehicles," in *2017 Resilience Week (RWS)*, 2017.
- [35] B. Kim, Y. Kashiba, S. Dai, and S. Shiraishi, "Testing autonomous vehicle software in the virtual prototyping environment," *Embedded Systems Letters*, vol. 9, no. 1, pp. 5–8, 2017.
- [36] B. Kim, A. Jarandikar, J. Shum, S. Shiraishi, and M. Yamaura, "The SMT-based automatic road network generation in vehicle simulation environment," in *International Conference on Embedded Software (EMSOFT)*. ACM, 2016, pp. 18:1–18:10.
- [37] C. Fan, B. Qi, and S. Mitra, "Data-driven formal reasoning and their applications in safety analysis of vehicle autonomy features," *IEEE Design Test*, vol. 35, no. 3, pp. 31–38, June 2018.
- [38] M. O'Kelly, H. Abbas, S. Gao, S. Shiraishi, S. Kato, and R. Mangharam, "APEX: Autonomous vehicle plan verification and execution," in *SAE World Congress*, 2016.
- [39] D. Cireşan, U. Meier, J. Masci, and J. Schmidhuber, "A committee of neural networks for traffic sign classification," in *Proceedings of the 2011 International Joint Conference on Neural Networks (IJCNN)*, 2011, pp. 1918–1921.
- [40] V. John, K. Yoneda, B. Qi, Z. Liu, and S. Mita, "Traffic light recognition in varying illumination using deep learning and saliency map," in *Proceedings of the 2014 IEEE 17th International Conference on Intelligent Transportation Systems (ITSC)*, 2014.
- [41] A. Angelova, A. Krizhevsky, and V. Vanhoucke, "Pedestrian detection with a large-field-of-view deep network," in *Proceedings of the IEEE International Conference on Robotics and Automation (ICRA '15)*, 2015.
- [42] S. Dutta, S. Jha, S. Sankaranarayanan, and A. Tiwari, "Learning and verification of feedback control systems using feedforward neural networks," in *Analysis and Design of Hybrid Systems*, 2018.
- [43] T. Dreossi, A. Donze, and S. A. Seshia, "Compositional falsification of cyber-physical systems with machine learning components," in *NASA Formal Methods (NFM)*, ser. LNCS, vol. 10227. Springer, 2017, pp. 357–372.
- [44] T. Dreossi, S. Jha, and S. A. Seshia, "Semantic adversarial deep learning," vol. 10981, pp. 3–26, 2018.
- [45] H. Abbas, M. O'Kelly, A. Rodionova, and R. Mangharam, "Safe at any speed: A simulation-based test harness for autonomous vehicles," in *7th International Workshop on Cyber-Physical Systems (CyPhy)*, 2017.
- [46] R. Alur, *Principles of Cyber-Physical Systems*. MIT Press, 2015.
- [47] H. Abbas, B. Hoxha, G. Fainekos, and K. Ueda, "Robustness-guided temporal logic testing and verification for stochastic cyber-physical systems," in *IEEE International Conference on CYBER Technology in Automation, Control, and Intelligent Systems*, 2014.
- [48] E. Bartocci, J. Deshmukh, A. Donzé, G. Fainekos, O. Maler, D. Nickovic, and S. Sankaranarayanan, "Specification-based monitoring of cyber-physical systems: A survey on theory, tools and applications," in *Lectures on Runtime Verification - Introductory and Advanced Topics*, ser. LNCS. Springer, 2018, vol. 10457, pp. 128–168.
- [49] G. E. Fainekos and G. J. Pappas, "Robustness of temporal logic specifications for continuous-time signals," *Theoretical Computer Science*, vol. 410, no. 42, pp. 4262–4291, 2009.
- [50] S. Boyd and L. Vandenberghe, *Convex Optimization*. Cambridge University Press, 2004.
- [51] T. A. Henzinger, P. W. Kopke, A. Puri, and P. Varaiya, "What's decidable about hybrid automata?" *J. Comput. Syst. Sci.*, vol. 57, no. 1, pp. 94–124, 1998.
- [52] B. Hoxha, H. Bach, H. Abbas, A. Dokhanchi, Y. Kobayashi, and G. Fainekos, "Towards formal specification visualization for testing and monitoring of cyber-physical systems," in *International Workshop on Design and Implementation of Formal Tools and Systems*, 2014.
- [53] J. Kapinski, J. V. Deshmukh, X. Jin, H. Ito, and K. Butts, "Simulation-based approaches for verification of embedded control systems: An overview of traditional and advanced modeling, testing, and verification techniques," *IEEE Control Systems Magazine*, vol. 36, no. 6, pp. 45–64, 2016.
- [54] D. R. Kuhn, R. N. Kacker, and Y. Lei, *Introduction to combinatorial testing*. CRC press, 2013.
- [55] S. Shalev-Shwartz, S. Shammah, and A. Shashua, "On a formal model of safe and scalable self-driving cars," arXiv:1708.06374v2, Tech. Rep., 2017.
- [56] R. Schubert, E. Richter, and G. Wanielik, "Comparison and evaluation of advanced motion models for vehicle tracking," in *2008 11th International Conference on Information Fusion*, June 2008, pp. 1–6.

- [57] D. J. LeBlanc, M. Gilbert, S. Stachowski, D. Blower, C. A. C. Flannagan, S. Karamihas, and W. T. B. andRini Sherony, "Advanced surrogate target development for evaluating pre-collision systems," in *23rd Enhanced Safety of Vehicles Conferences*, 2013.
- [58] A. Petrovskaya and S. Thrun, "Model based vehicle detection and tracking for autonomous urban driving," *Autonomous Robots*, vol. 26, no. 2, pp. 123–139, Apr 2009.
- [59] B. Grabe, T. Ike, and M. Hötter, "Evidence based evaluation method for grid-based environmental representation," in *FUSION*. IEEE, 2009, pp. 1234–1240.
- [60] "ISO Functional Safety," https://en.wikipedia.org/wiki/ISO_26262, accessed: 2018-09-11.
- [61] M. Elbanhawi, M. Simic, and R. Jazar, "In the passenger seat: Investigating ride comfort measures in autonomous cars," *IEEE Intelligent Transportation Systems Magazine*, vol. 7, no. 3, pp. 4–17, Fall 2015.
- [62] P. Green, "Motion sickness and concerns for self-driving vehicles: A literature review," <http://umich.edu/simdriving/publications/Motion-Sickness--Report-061616pg-sent.pdf>, accessed: 2018-09-11.
- [63] A. Dokhanchi, B. Hoxha, C. E. Tuncali, and G. Fainekos, "An efficient algorithm for monitoring practical tptl specifications," in *Formal Methods and Models for System Design (MEMOCODE), 2016 ACM/IEEE International Conference on*. IEEE, 2016, pp. 184–193.
- [64] O. Michel, "Cyberbotics ltd. Webots: professional mobile robot simulation," *International Journal of Advanced Robotic Systems*, vol. 1, no. 1, p. 5, 2004.
- [65] M. Abadi, A. Agarwal, P. Barham, E. Brevdo, Z. Chen, C. Citro, G. S. Corrado, A. Davis, J. Dean, M. Devin, *et al.*, "Tensorflow: Large-scale machine learning on heterogeneous distributed systems," *arXiv preprint:1603.04467*, 2016.
- [66] M. Ester, H.-P. Kriegel, J. Sander, X. Xu, *et al.*, "A density-based algorithm for discovering clusters in large spatial databases with noise." in *Kdd*, vol. 96, no. 34, 1996, pp. 226–231.
- [67] E. Schubert, J. Sander, M. Ester, H. P. Kriegel, and X. Xu, "DbSCAN revisited, revisited: why and how you should (still) use dbSCAN," *ACM Transactions on Database Systems (TODS)*, vol. 42, no. 3, p. 19, 2017.
- [68] E. A. Wan and R. Van Der Merwe, "The unscented kalman filter for nonlinear estimation," in *Proceedings of the IEEE 2000 Adaptive Systems for Signal Processing, Communications, and Control Symposium (Cat. No. 00EX373)*. Ieee, 2000, pp. 153–158.
- [69] A. Donze, "Breach, a toolbox for verification and parameter synthesis of hybrid systems," in *Computer Aided Verification (CAV)*, ser. LNCS, vol. 6174. Springer, 2010, pp. 167–170.
- [70] T. Dreossi, S. Ghosh, X. Yue, K. Keutzer, A. L. Sangiovanni-Vincentelli, and S. A. Seshia, "Counterexample-guided data augmentation," in *Proceedings of the Twenty-Seventh International Joint Conference on Artificial Intelligence, IJCAI*, 2018, pp. 2071–2078.
- [71] A. Dosovitskiy, G. Ros, F. Codevilla, A. Lopez, and V. Koltun, "CARLA: An open urban driving simulator," in *Proceedings of the 1st Annual Conference on Robot Learning*, 2017, pp. 1–16.
- [72] B. Hoxha, N. Mavridis, and G. Fainekos, "VISPEC: A graphical tool for elicitation of MTL requirements," in *Proceedings of the IEEE/RSJ International Conference on Intelligent Robots and Systems*, 2015.